

Indirect Associations in Learning Semantic and Syntactic Lexical Relationships

M. A. Kelly^{a,b,*}, Moojan Ghafurian^{a,c}, Robert L. West^d, David Reitter^{a,e}

^a*The Pennsylvania State University, University Park, PA*

^b*Bucknell University, Lewisburg, PA*

^c*University of Waterloo, Waterloo, ON, Canada*

^d*Carleton University, Ottawa, ON, Canada*

^e*Google Research, New York City, NY*

Abstract

Computational models of distributional semantics (a.k.a. word embeddings) represent a word’s meaning in terms of its relationships with all other words. We examine what grammatical information is encoded in distributional models and investigate the role of indirect associations. Distributional models are sensitive to associations between words at one degree of separation, such as ‘tiger’ and ‘stripes’, or two degrees of separation, such as ‘soar’ and ‘fly’. By recursively adding higher levels of representations to a computational, holographic model of semantic memory, we construct a distributional model sensitive to associations between words at arbitrary degrees of separation. We find that word associations at four degrees of separation increase the similarity assigned by the model to English words that share part-of-speech or syntactic type. Word associations at four degrees of separation also improve the ability of the model to construct grammatical English sentences. Our model proposes that human memory uses indirect associations to learn part-of-speech and that the basic associative mechanisms of memory and learning support knowledge of both semantics and grammatical structure.

Keywords: distributional semantics, semantic memory, word embeddings,

*Corresponding author at: Department of Computer Science, Bucknell University, Dana 336, Lewisburg, PA 17837, USA.

Email addresses: `m.alex.kelly@bucknell.edu` (M. A. Kelly), `moojan@uwaterloo.ca` (Moojan Ghafurian), `robert.west@carleton.ca` (Robert L. West), `reitter@google.com` (David Reitter)

1. Introduction

Syntax (how words are put together) and semantics (what words mean) have traditionally been understood as arising from distinct cognitive processes. The distinction between syntax and semantics was famously illustrated by Chomsky (1956) with the example “Colorless green ideas sleep furiously”, a sentence that is grammatical but meaningless.

But can syntax and semantics be understood as arising from a unitary cognitive process? Predictive neural language models (e.g., Ororbia II et al., 2017) appear to be sensitive to both syntax and semantics. Recurrent neural networks are able to make judgements about subject-verb agreement in nonsensical sentences such as “Colorless green ideas sleep furiously” without needing to rely on part-of-speech tagging or other syntactic markers (Gulordava et al., 2018). However, due to the “black box” nature of neural network models, it is difficult to say exactly what information is being exploited by the networks to make decisions about syntax.

Even though the nonsense sentence “Colorless green ideas sleep furiously” has a set of word transitions that do not appear in English language corpora, the sentence has a very common English construction: *adjective, adjective, noun, verb, adverb*. How do humans learn that, at an abstract level, the sentence is structurally similar to many other sentences in their life experience?

Jenkins (1964, 1965) and Jenkins and Palermo (1964) hypothesize that knowledge of the syntactic structure of language depends on indirect or *mediated* associations. More specifically, part-of-speech, or the knowledge that nouns can be substituted for other nouns and verbs for other verbs, and so on, depends on learning equivalence classes through mediated association. Although Jenkins (1968, 1974) ultimately abandoned the paradigm of understanding language and memory in terms of associations and equivalence classes altogether, more recent studies with children have found that exploiting equivalence classes is a powerful pedagogical technique for rapidly expanding a learner’s language abilities (Sidman, 2009).

To explore the hypothesis that learning the part-of-speech of words is based on a capacity for indirect or mediated association, we propose a “deep” distributional semantics model, the *Hierarchical Holographic Model* (HHM). HHM consists of a stack of holographic vector models that feed one into the

next, which allows HHM to detect arbitrarily indirect associations between words. HHM is based on BEAGLE (Jones et al., 2006; Jones and Mewhort, 2007), one of the few distributional semantics models sensitive to the order of words in sentences, a critical part of English syntax.

Holographic models of human memory have a long history (Murdock, 1982; Pribram, 1969) and have been applied to a wide range of paradigms (e.g., Eliasmith, 2013; Franklin and Mewhort, 2015; Jamieson and Mewhort, 2011). Holographic vectors allow for easy implementation of a recursive model capable of learning arbitrarily indirect associations. Our approach can be understood as an extension of Jenkins and Palermo (1964)’s work, though instead of using artificial grammar experiments, we use a computational approach applied to an English-language corpus.

In what follows, we provide theoretical background on the Hierarchical Holographic Model and then evaluate the model. We give a proof-of-concept demonstration of HHM on a small artificial dataset and then train HHM on an English-language corpus. We analyze the relationship between the representations produced by the higher levels of HHM and part-of-speech (e.g., nouns, adjectives, adverbs, etc.) and the syntactic types proposed by Combinatory Categorical Grammar (CCG; Steedman and Baldridge, 2011). We show that HHM’s representations can be used to order words into grammatical sentences and we test HHM on the sentence “Colorless green ideas sleep furiously”. HHM is an account of the mental lexicon based on a general-purpose computational model of human memory. HHM demonstrates how a single system can incorporate knowledge of both how a word is used (i.e., part-of-speech) and what a word means (i.e., distributional semantics).

2. Theory

In this section, we describe the BEAGLE model of distributional semantics (Jones and Mewhort, 2007), based on the holographic model of memory (Plate, 1995). We propose the Hierarchical Holographic Model (HHM). HHM is a recursively constructed variant of BEAGLE capable of detecting arbitrarily high orders of association. We then define *orders of association* as a measure of the relationship between a pair of words in memory.

2.1. The BEAGLE Model

The BEAGLE model (Jones and Mewhort, 2007) belongs to the family of distributional semantics models, also known as *word embeddings*. Distri-

butional models include Latent Semantic Analysis (Landauer and Dumais, 1997), the Hyperspace Analogue to Language (Burgess and Lund, 1997), the Topics Model (Griffiths et al., 2007), *word2vec* (Mikolov et al., 2013), GloVe (Pennington et al., 2014), as well as word embeddings extracted from neural language models such as BERT (Devlin et al., 2019). Distributional models use the word co-occurrence statistics of a large corpus to construct high-dimensional vectors that represent the meanings of words. Each vector can be understood as a point in a high-dimensional space and distance in the space serves as a measure of similarity in meaning. Words that are closer together have more similar meanings. Such a space, where distance measures similarity in meaning, is referred to as a *semantic space*.

In BEAGLE, each word is represented by two vectors: an *environment vector* that represents the percept of a word (i.e., the word’s perceptual features) and a *memory vector* that represents the concept of a word (i.e., the word’s meaning and associations).

An environment vector (denoted by **e**) stands for what a word looks like in writing or sounds like when spoken. For simplicity, we do not simulate the visual or auditory features of words (but see Cox et al., 2011, for a version of BEAGLE that does simulate features). Instead, we generate the environment vectors using random values, as in Jones and Mewhort (2007). Thus, in our simulations, words with similar morphology (e.g., *walk* and *walked*) have dissimilar environment vectors, such that the model needs to learn from the corpus that the two words are related.

Environment vectors are generated by randomly sampling values from a Gaussian distribution with a mean of zero and a variance of $1/d$, where d is the dimensionality. Individually, the dimensions of the vectors have no inherent meaning: they do not stand for specific words or features. A word is represented as a pattern of values across all dimensions. The number of dimensions, d , determines the fidelity with which BEAGLE stores word co-occurrence information, such that smaller d yields poorer encoding.

Memory vectors (denoted by **m**) represent the associations a word has with other words. As the model reads the corpus, memory vectors are continuously updated. For example, the words *walk* and *walked* are represented by dissimilar, randomly-generated environment vectors. But, because the words are used in similar ways, *walk* and *walked* develop highly similar memory vectors. That said, the two memory vectors will not be identical, as *walked* is more likely to appear in contexts with other past-tense verbs and *walk* with other present-tense verbs (e.g., “I *walked* to the store and *bought* bread” vs.

108 “I *walk* to the store and *buy* bread.”).

109 BEAGLE stores two kinds of information in a memory vector: *context*
110 and *order*. The *context* information for a target word in a sentence is the sum
111 of the environment vectors of the other words in the sentence. Conversely,
112 the *order* information for a word in a sentence is a sum of sequences of words
113 that include the target word. A sequence of words is represented by a vector
114 that is a *convolution* of the environment vectors of the words in the sequence.

115 2.1.1. Order information

116 The memory vectors are termed *holographic* because they use circular
117 convolution to compactly encode associations between words (Plate, 1995).
118 According to holographic theories of memory (Eliasmith, 2013; Murdock,
119 1982; Pribram, 1969), patterns of neural activity in the brain interfere to
120 create new associations in a manner mathematically analogous to how light
121 waves interfere to create a hologram (Gabor, 1969). Given two patterns of
122 neural activity represented as vectors, the interference pattern produced by
123 the association of the two is computed as the *convolution* of the vectors.

124 To compute the *order* information for a target word, a sum of n -grams
125 is added to the target word’s memory vector. The n -grams are at minimum
126 bigrams consisting of the target word and the word immediately preceding or
127 following. The n -grams also have a maximum size that can be set. Jones and
128 Mewhort (2007) use a maximum of 7-grams. We experiment with maximum
129 n -gram sizes ranging from 5-grams to the full length of the sentence.

130 For example, given the sentence, “eagles soar over trees”, BEAGLE up-
131 dates the memory vectors for each word in the sentence: *eagles*, *soar*, *over*,
132 and *trees*. For *soar*, the following n -grams are added into the memory vector
133 \mathbf{m}_{soar} : the bigrams “eagles soar” and “soar over”, the trigrams “eagles soar
134 over” and “soar over trees”, and the tetragram “eagles soar over trees”.

135 Each n -gram is constructed as a convolution of the environment vectors
136 of the constituent words, except for the target word, which is represented by
137 the placeholder vector (denoted by Φ). The placeholder vector is randomly
138 generated and serves as a universal retrieval cue. With the placeholder sub-
139 stituted for the target word, each n -gram can be understood as a question to
140 which the target word is the answer. So, rather than adding a representation
141 of “eagles soar over” into \mathbf{m}_{soar} , we instead add “eagles ? over”, i.e., “What
142 was the word that appeared between *eagles* and *over*?”. Each memory vector
143 can be understood as the sum of all questions to which that memory vector’s
144 word is an appropriate answer.

145 Given “eagles soar over trees”, we add “eagles ?”, “? over”, “eagles ?
 146 over”, “? over trees”, and “eagles ? over trees” to \mathbf{m}_{soar} as follows:

$$\begin{aligned}\mathbf{m}_{\text{soar},t+1} = & \mathbf{m}_{\text{soar},t} + (\mathbf{P}_{\text{before}} \mathbf{e}_{\text{eagles}}) * \Phi + \\ & (\mathbf{P}_{\text{before}} \Phi) * \mathbf{e}_{\text{over}} + \\ & (\mathbf{P}_{\text{before}} ((\mathbf{P}_{\text{before}} \mathbf{e}_{\text{eagles}}) * \Phi)) * \mathbf{e}_{\text{over}} + \\ & (\mathbf{P}_{\text{before}} ((\mathbf{P}_{\text{before}} \Phi) * \mathbf{e}_{\text{over}})) * \mathbf{e}_{\text{trees}} + \\ & (\mathbf{P}_{\text{before}} ((\mathbf{P}_{\text{before}} ((\mathbf{P}_{\text{before}} \mathbf{e}_{\text{eagles}}) * \Phi)) * \mathbf{e}_{\text{over}})) * \mathbf{e}_{\text{trees}}\end{aligned}$$

147 where $*$ is circular convolution, t is the current time step, all vectors \mathbf{m} , \mathbf{e} , and
 148 Φ have d dimensions, and $\mathbf{P}_{\text{before}}$ is a permutation matrix used to indicate
 149 that a word occurred earlier in the sequence (see Appendix for discussion).
 150 $\mathbf{P}_{\text{before}}$ is made by randomly reordering the rows of the $d \times d$ identity matrix.
 151 Multiplying a vector \mathbf{v} by $\mathbf{P}_{\text{before}}$ results in the permuted vector $\mathbf{P}_{\text{before}} \mathbf{v}$.

152 2.1.2. Context information

153 *Context* information is a sum of environment vectors. For example, the
 154 context information for \mathbf{m}_{soar} and the sentence “Eagles soar over trees” is:

$$\mathbf{m}_{\text{soar},t+1} = \mathbf{m}_{\text{soar},t} + \mathbf{e}_{\text{eagles}} + \mathbf{e}_{\text{over}} + \mathbf{e}_{\text{trees}} \quad (1)$$

155 For the purposes of the simulations reported in this paper, we only use the
 156 *order* information and exclude the *context* information, as we found little
 157 benefit to including *context* information in the word ordering task that we
 158 use to evaluate the models.

159 2.1.3. Applications of BEAGLE

160 BEAGLE can model semantic priming (Jones et al., 2006), the pattern
 161 of semantic memory deficits in Alzheimer’s disease (Johns et al., 2013), as
 162 well as basic memory phenomena, such as release from proactive interference
 163 (Mewhort et al., 2018).

164 While BEAGLE is a model of the mental lexicon, Dynamically Structured
 165 Holographic Memory (Rutledge-Taylor et al., 2014) is a variant of BEA-
 166 GLE applied to non-linguistic memory and learning tasks, such as learn-
 167 ing sequences of actions for strategic game play. Kelly et al. (2015) and
 168 Kelly and Reitter (2017) propose another BEAGLE variant, Holographic
 169 Declarative Memory, that learns sets of property-value pairs (e.g., *colour:red*

170 *shape:octagon type:sign*) of the kind used by the ACT-R cognitive architec-
171 ture (Anderson, 2009), showing that BEAGLE’s algorithm can be applied to
172 any problem domain that can be expressed in discrete symbols.

173 The Hierarchical Holographic Model (HHM) can, like BEAGLE, be ap-
174 plied to a wide range of problem domains. While we evaluate HHM in this
175 paper in terms of its ability to account for properties of natural language,
176 HHM is intended as a general model of learning and memory.

177 2.2. Hierarchical Holographic Model

178 The Hierarchical Holographic Model (HHM) is a series of BEAGLE-like
179 models, such that the memory vectors of one model serve as the environ-
180 ment vectors for the next model. Level 1 is a standard BEAGLE model
181 with randomly generated environment vectors, except that we only use or-
182 der information to construct the memory vectors. Level 2 and higher are
183 order-only BEAGLE models where the environment vectors are the memory
184 vectors of the previous level. Once Level 1 has been run on a corpus, Level 2
185 is initialized with Level 1’s memory vectors as its environment vectors. Then
186 Level 2 is run on the corpus to generate a new set of memory vectors, which
187 in turn are used as the environment vectors for the next level, and so on, to
188 generate as many levels of representations as desired.

189 To use the memory vectors of a previous level as the environment vectors
190 for the next, one must normalize and randomly permute the vectors. Vectors
191 are normalized to unit Euclidean length to ensure that each word is equally
192 weighted at the next level. Without normalization, high-frequency words
193 would disproportionately dominate the representations at the next level.

194 Permutation is necessary to protect the information encoded at one level
195 from information encoded at the next level (Gayler, 2003). Without using
196 permutation, the different levels of information become confounded and de-
197 structively interfere with each other (Kelly et al., 2013). The destructive
198 interference arises because convolution distributes over addition. If we con-
199 volve a memory vector with another vector, that vector will distribute across
200 all of the component n -gram vectors that are summed into the memory vec-
201 tor. If the other vector is also a memory vector, all of its n -grams will
202 distribute across all of the memory vector’s n -grams to create a multitude of
203 spurious n -gram representations.

204 Thus, to transform memory vectors to environment vectors, the elements
205 of all memory vectors are re-ordered according to a randomly generated per-

206 mutation, $\mathbf{P}_{\text{group}}$. For level $l + 1$, and all words i , the environment vectors
 207 for that level are:

$$\mathbf{e}_{l+1,i} = \mathbf{P}_{\text{group}}\left(\frac{\mathbf{m}_{l,i}}{\sqrt{\mathbf{m}_{l,i} \bullet \mathbf{m}_{l,i}}}\right) \quad (2)$$

208 where \mathbf{e} and \mathbf{m} are environment and memory vectors and \bullet is the dot product.

209 The levels in HHM can be understood as the products of memory re-
 210 consolidation, the process of revisiting experiences and recording new infor-
 211 mation about those experiences. The different levels of representation are
 212 stored separately from each other in the model for the purpose of examining
 213 the differential effects of representations that encode lower and higher or-
 214 ders of associations. The different levels are not necessarily separate memory
 215 systems, but instead could constitute different kinds of knowledge within a
 216 single memory system.

217 2.3. Orders of Association

218 Saussure (1916) defines two types of relationships between words: *paradig-*
 219 *matic* and *syntagmatic*. *Syntagmatic* describes a relationship a word has with
 220 other words that surround it. *Paradigmatic* describes a relationship in which
 221 a pair of words can be substituted for each other.

222 Grefenstette (1994) defines first-order, second-order, and third-order affini-
 223 ties between words and notes that computational language models are typi-
 224 cally sensitive to either first-order (topic) or second-order (synonymy) affini-
 225 ties. Grefenstette defines third-order affinities as semantic groupings among
 226 similar words, which can be discovered using cluster analysis techniques.

227 We define the term *order of association* as a measure of the degree of
 228 separation of two words in an agent’s language experience. Imagine a graph
 229 where each word in the lexicon is a node connected to other words. *Order*
 230 *of association* is the length of a path between two words in the graph. The
 231 *strength* of that order of association is the number of paths of that length
 232 between the two words.

233 A pair of words are connected once for each time they have occurred in the
 234 same context. In human cognition, the context is defined by the associations
 235 in mind at the time of encoding. Ideally, we would use a model of memory
 236 to determine when words are or are not in the same context (see §5.2 for
 237 discussion). However, for simplicity, we use a context that is a window of
 238 five or more words to the left and right of the target word.

239 *First-order association* describes two words that appear together. In the
 240 sentence “eagles soar over trees”, the words *eagles* and *trees* have first-order
 241 association. Words with strong first-order association (i.e., frequently appear
 242 together) are often related in topic (i.e., have a *syntagmatic* relationship),
 243 such as the words *tiger* and *stripes*.

244 *Second-order association* describes two words that appear with the same
 245 words. Given “airplanes soar through skies” and “airplanes fly through
 246 skies”, *soar* and *fly* have second-order association. Words with strong second-
 247 order association are often synonyms (i.e., have a *paradigmatic* relationship).

248 *Third-order association* is a first-order association plus a second-order as-
 249 sociation (i.e., a paradigmatic relationship plus a syntagmatic relationship).
 250 For example, *tiger* and *stripes* have a first-order association and *lion* and *tiger*
 251 have a second-order association. Thus, *lion* and *stripes* have a third-order
 252 association mediated by *tiger*.

253 Statistical smoothing algorithms use third-order associations to estimate
 254 the acceptability of novel bigrams (Pereira, 2000; Roberts and Chater, 2008).
 255 For example, *unsightly bumbershoot* is a perfectly acceptable adjective-noun
 256 pair, but is unlikely to appear in a corpus that doesn’t include this paper. But
 257 an *unsightly bumbershoot* is very similar to an *unsightly umbrella*. The third-
 258 order association between *unsightly* and *bumbershoot* mediated by *umbrella*
 259 can be used to judge that *unsightly bumbershoot* is an acceptable bigram.

260 *Fourth-order association* describes two words that appear with words that
 261 appear with the same words. A fourth-order association is two second-order
 262 (or paradigmatic) associations added together.

263 The sentences in Table 1 provide an artificial example of a fourth-order
 264 association. Words with fourth-order association are indicated in **bold** and
 265 words with second-order association are indicated in *italics*. The word pairs
 266 *soar* and *fly*, *over* and *above*, and *trees* and *forest* each have second-order
 267 associations. Given only the sentences in Table 1, the words *eagles* and *birds*
 268 do not have first-, second-, or third-order association, but do have fourth-
 269 order. The web of associations between the words in Table 1’s sentences is
 270 illustrated in Figure 1.

271 Table 1 is an artificial example. In natural language, *eagles* and *birds*
 272 have strong second-order association (i.e., are highly synonymous). Fourth-
 273 order association indicates that two words can be substituted for each other,
 274 but at a more abstract level than second-order association (synonymy). We
 275 hypothesize that word pairs that have strong fourth-order association, but do
 276 not have first- or second-order association, are words unrelated in meaning

Table 1: Example of a fourth order association between *eagles* and *birds*.

Sentences	
eagles <i>soar over trees</i>	birds <i>fly above forest</i>
airplanes <i>soar</i> through skies	airplanes <i>fly</i> through skies
dishes are <i>over</i> plates	dishes are <i>above</i> plates
squirrels live in <i>trees</i>	squirrels live in <i>forest</i>
cars drive on streets	

277 but are grammatically acceptable to substitute for each other. We expect
 278 that words with fourth-order association are likely to share part-of-speech
 279 or syntactic type (e.g., *focused* and *emerging* can both be used as a verb or
 280 adjective, see Table 2). We explore this hypothesis in Sections 3.3 and 3.4.

281 *Fifth-order and higher* associations can be obtained by abstracting indef-
 282 initely. Eventually, all words are related to all other words in the language.

283 *Even-numbered* associations are paradigmatic or *super-paradigmatic* rela-
 284 tionships that indicate a semantically valid or, we hypothesize, syntactically
 285 valid substitution.

286 *Odd-numbered* associations are syntagmatic or *super-syntagmatic* rela-
 287 tionships describing the association between a word and other words that
 288 could appear either with the word directly (first-order) or with another word
 289 like it (third-order, fifth-, etc.).

290 *No association* describes a pair of words that have no path between them
 291 of any length. For an agent that knows only the nine sentences in Table 1,
 292 the words *car* and *eagle* have no association. In real language data, two
 293 words will only have no association if they belong to two different languages
 294 (e.g., the words *goyang-i* from Korean and *borroka* from Basque have no
 295 association with each other).

296 In our description of *orders of association* we have glossed over the ques-
 297 tion of the distinct nature of syntagmatic versus paradigmatic associations.
 298 For two words to have a syntagmatic association, it is sufficient for the words
 299 to co-occur in any way. Conversely, for paradigmatic associations, the two
 300 words should be interchangeable for each other, which is contingent on posi-
 301 tion in the sentence or phrase.

302 HHM, as implemented in this paper, is specifically a model of super-
 303 paradigmatic associations between words. Examining super-syntagmatic as-

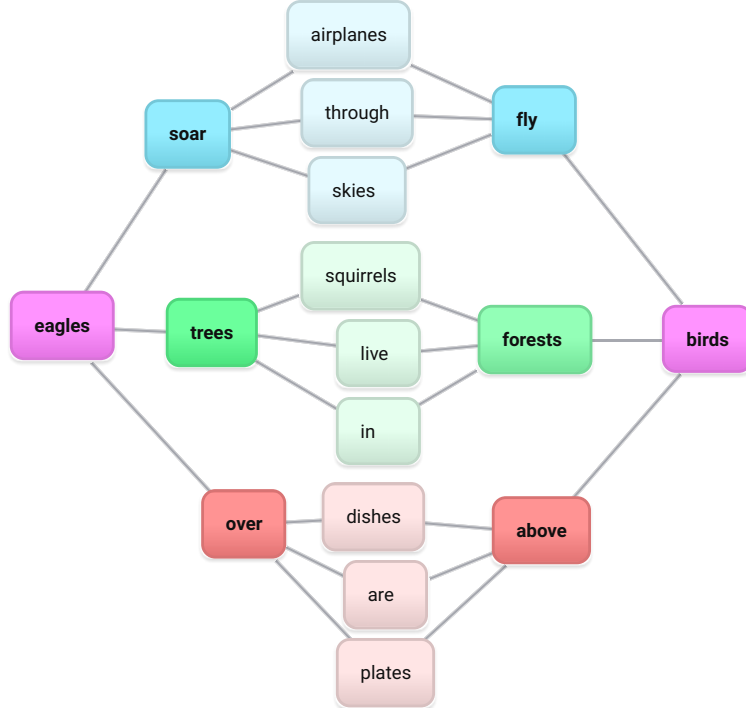


Figure 1: Web of associations between words in Table 1.

sociations is beyond the scope of our work, as our interest is in part-of-speech and syntactic type relationships, which are valid substitution relationships, rather than co-occurrence (or syntagmatic) relationships. For the purposes of this paper, we only use *order* vectors in HHM. However, we have found that odd-numbered orders of association are captured by recursively constructing levels of representation using *context* vectors.

To define orders of association, we have described the lexicon as a connected graph. This graph is not explicitly represented by HHM. HHM defines a semantic space rather than a graph. Words close together at Level 1 of HHM have strong second-order association, Level 2 represents fourth-order associations, Level 3 represents sixth-order associations, and so on.

Note that *order of association* in a language is distinct from *orders of approximation* to a language. *Orders of approximation* is a measure of how closely a probability model approximates a language as measured by the number of words that are taken into account when predicting the next word

in a sequence (Shannon, 1951). Depending on the size of the HHM context window, we use up to 5, 10, or k preceding words to predict a word as well as up to 5, 10, or k of the succeeding words, where k is the length of the sentence. As such, HHM could be described as a 5th, 10th, or k th order approximation to English. Independent of this parameter is the order of association. In this paper, we explore using up to eighth-order associations. Order of approximation and association interact, such that higher orders of approximation (i.e., larger context windows) are more useful in a model sensitive to higher orders of association.

3. Simulations and Experiments

We test two hypotheses:

1. Level 2 (fourth-order associations) or higher levels of the Hierarchical Holographic Model (HHM) significantly outperform Level 1 (second-order associations) on tests of correlates of syntactic knowledge.
2. Whereas second-order associations are semantic in character, fourth-order associations or higher provide knowledge that is primarily part-of-speech or a word’s syntactic type.

We contrast the two hypotheses with two alternatives:

1. Fourth-order associations or higher do not improve performance on tests of correlates of syntactic knowledge.
2. Fourth-order associations or higher merely provide additional lexical semantic knowledge, such that given more data, a model sensitive only to second-order associations would discover the same word relationships.

To test these hypotheses, we begin by validating HHM as a model of orders of association. We show that HHM works as intended and is able to detect fourth-order associations in a small artificial data set (Section 3.1).

To demonstrate that higher-order associations are lexical syntactic in character, we investigate the relationship between higher-order associations and part-of-speech (Experiment 1).

However, part-of-speech provides only a coarse-grained analysis of the types of words in English. Conversely, Combinatory Categorical Grammar (CCG; Steedman and Baldrige, 2011) postulates hundreds of different word types. In CCG, a word type captures what types of phrases the word may

352 combine with to the left or to the right (and the associated semantic op-
353 erations). Thus grammatical information is stored along with the word in
354 the lexicon, providing fine-grained information about how each word is used.
355 The theory proposes a very limited set of syntactic and semantic operations
356 in parsing and sentence production that is parameterized for the specific
357 language. CCG is a broad-coverage formalism that allows us to study the
358 granularity of grammatical information that might be represented in the vec-
359 tors generated by higher-order associations (Experiment 2).

360 While comparisons between HHM, part-of-speech, and CCG types are
361 illuminating, part-of-speech and CCG are theories of language, not language
362 itself. To evaluate the role of higher-order associations in producing gram-
363 matical sentences, we situate HHM’s word representations in a simple exem-
364 plar model that operates on sentences. We use a word ordering task where the
365 exemplar model must order a given set of words into a grammatical sentence.
366 By varying the level of HHM used by the exemplar model, we investigate the
367 effect of higher-orders of association and n -gram size on the ability of the
368 model to find the grammatical ordering of the words (Experiment 3).

369 Chomsky (1956) famously gave “Colorless green ideas sleep furiously”
370 as an example of a sentence that is grammatical but meaningless. If the
371 sentence is truly meaningless we would expect that second-order (semantic)
372 associations would be insufficient for finding the grammatical ordering of
373 the words *colorless*, *furiously*, *green*, *ideas*, and *sleep*. However, if fourth-
374 order associations are syntactic in character, we should expect to find that
375 the exemplar model can find the grammatical ordering of the words using
376 representations from Level 2 of HHM (Experiment 4).

377 Through these simulations and experiments, we seek to demonstrate the
378 validity of HHM as a model, HHM’s relationship to established theories of
379 syntax, and the role of higher-order associations in constructing grammatical
380 sentences. Code for running HHM¹² and the exemplar model is available
381 online, along with data and figures³.

382 3.1. *Small Example on Artificial Data*

383 Here we show that HHM is able to detect higher-order associations as
384 intended. For the purposes of providing a clear illustration of the behavior

¹<https://github.com/ecphory/BEAGLE-HHM>

²<https://github.com/moojan/Python-BEAGLE-HHM>

³<https://github.com/ecphory/Indirect-Associations>

385 of the model, we use a small artificial data set that provides a clean example
 386 of first-, second-, and fourth-order associations. The data set consists exclu-
 387 sively of the sentences in Table 1. This is merely a toy example, but useful
 388 for demonstrating how the model works. This example has been designed
 389 such that the word pairs *soar* and *fly*, *over* and *above*, and *trees* and *forest*,
 390 have second-order associations, whereas the word pair *eagles* and *birds*, have
 391 a fourth-order association.

392 HHM was run with 1024 dimensional vectors and three levels of repre-
 393 sentations. In the nine sentences of this example, there are 21 unique words,
 394 and thus 210 unique pairs of words. We can characterize the behavior of
 395 HHM by how the word pairs change in similarity across levels.

396 Figure 2 shows cosine similarity between the word pairs as a function of
 397 level of representation in HHM. Of the 210 word pairs, we graph the 24 word
 398 pairs that have non-negative similarity by Level 3. Of those 24 pairs, we
 399 label and rank the 10 pairs with the most similarity, from *over above* (cosine
 400 = 0.70 at Level 3) to *over in* (cosine = 0.20 at Level 3). Word pairs with
 401 fourth-order association are in **bold** and word pairs with strong second-order
 402 association are in *italics*.

403 The memory vectors for words with second-order association are close on
 404 Level 1 (e.g., *soar* and *fly*, cosine = 0.51) and closer by Level 3 (cosine =
 405 0.67). The words *eagle* and *bird*, which have only fourth-order association,
 406 are unrelated on Level 1 (cosine = -0.01) but are the fifth most similar word
 407 pair by Level 3 (cosine = 0.33).

408 The results provide a simple example of the effect of the higher levels.
 409 Each memory vector at Level 1 is constructed as a sum of convolutions of
 410 environment vectors. As such, the memory vectors at Level 1 encode first-
 411 order associations with respect to the environment vectors, measuring the
 412 frequency with which each word co-occurs with other words and sequences
 413 of words. The cosines between memory vectors are a measure of second-
 414 order association, the degree to which the two words co-occur with the same
 415 words. The algorithm that produces Level 1 transforms data that captures
 416 first-order association (co-occurrence) into data that captures second-order
 417 associations. The algorithm is a step, and by repeating it to produce higher
 418 levels, we can build a staircase.

419 Level 1 of the model cannot detect associations higher than second-order.
 420 A pair of words with third-order association or higher, but not first or second,
 421 do not appear together in the same sentence and do not co-occur with the
 422 same words. As such, the memory vectors for a pair of words with only third-

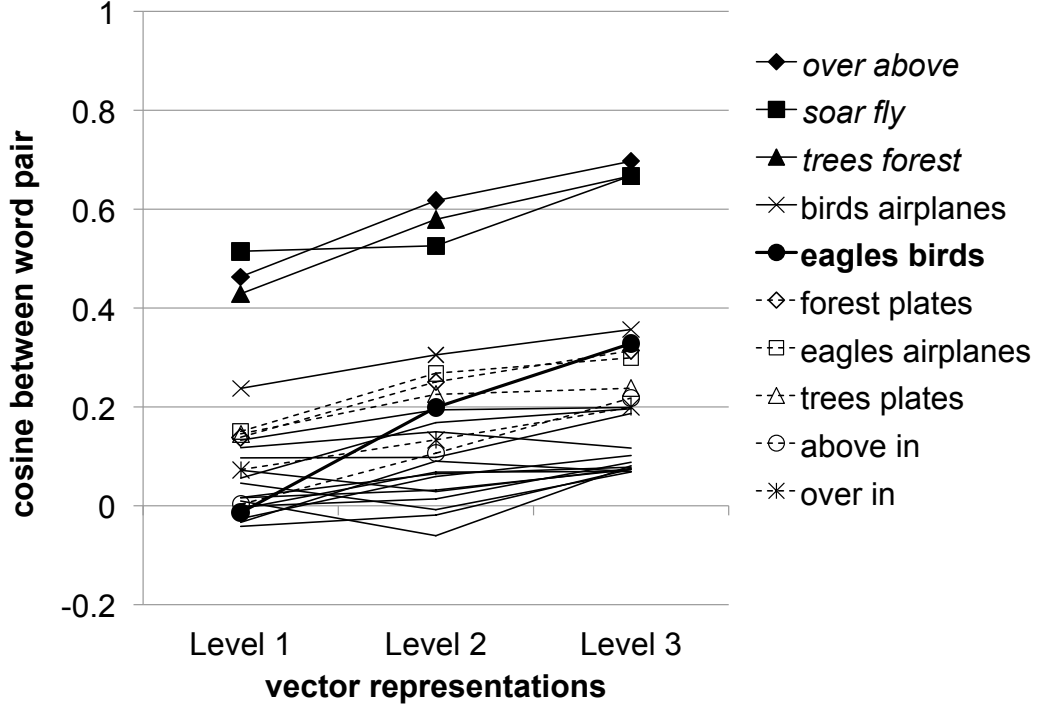


Figure 2: Cosines between word pairs across levels.

order or higher association will be constructed from disjoint sets of vectors. At Level 1, $\mathbf{m}_{1,\text{eagles}}$ is a sum of convolutions of $\mathbf{e}_{1,\text{soar}}$, $\mathbf{e}_{1,\text{over}}$, $\mathbf{e}_{1,\text{forest}}$, whereas $\mathbf{m}_{1,\text{birds}}$ is a sum of convolutions of $\mathbf{e}_{1,\text{fly}}$, $\mathbf{e}_{1,\text{above}}$, $\mathbf{e}_{1,\text{trees}}$. As Level 1 environment vectors are approximately orthogonal, the memory vectors constructed from them will also be approximately orthogonal. As a result, $\mathbf{m}_{1,\text{eagles}}$ and $\mathbf{m}_{1,\text{birds}}$ are approximately orthogonal (cosine = -0.01).

But at higher levels, the environment vectors are no longer orthogonal because the environment vectors for Level 2 are the memory vectors for Level 1. As a result, $\mathbf{e}_{2,\text{soar}}$ is similar to $\mathbf{e}_{2,\text{fly}}$ (cosine = 0.51), $\mathbf{e}_{2,\text{over}}$ is similar to $\mathbf{e}_{2,\text{above}}$ (cosine = 0.46), and $\mathbf{e}_{2,\text{forest}}$ is similar to $\mathbf{e}_{2,\text{trees}}$ (cosine = 0.43). Even though $\mathbf{m}_{2,\text{eagles}}$ and $\mathbf{m}_{2,\text{birds}}$ are still constructed from disjoint sets of environment vectors, because the vectors that they are constructed from are similar, $\mathbf{m}_{2,\text{eagles}}$ and $\mathbf{m}_{2,\text{birds}}$ are somewhat similar (cosine = 0.20).

Because the Level 2 environment vectors are more similar to each other than the Level 1 environment vectors, the memory vectors for the pairs *soar*

and *fly*, *above* and *over*, and *forest* and *trees* are also more similar at Level 2 than at Level 1 (see Figure 2). As a result, the Level 3 environment vectors for the three word pairs will be more similar at Level 3 than Level 2, which drives up the similarity between *eagles* and *birds* (cosine = 0.33). Eventually, at even higher levels, each pair *soar* and *fly*, *above* and *over*, and *forest* and *trees* will converge approximately to a point (cosine ≈ 1), causing *eagles* and *birds* to converge as well. The similarity between *eagles* and *birds* will never exceed the similarity between the three words pairs that the fourth-order association is contingent upon because it is the strengthening of those second-order associations that drives the strength of the fourth-order association.

3.2. Training the Model

We train HHM on the Novels Corpus from Johns et al. (2016b) with 10 238 600 sentences, 145 393 172 words, and 39 076 unique words. HHM reads the corpus one sentence at a time. Within each sentence, HHM uses a moving window centered on a target word. Within the window, all n -grams that include the target word, from bigrams up to n -grams of window width, are encoded as convolutions of environment vectors and summed into the target word’s memory vector. We use 1024 dimensional vectors and four levels of representations, where Level 1 is sensitive to second-order associations, Level 2 to fourth-order, Level 3 to sixth-order, and Level 4 to eighth-order.

At each level of HHM, we experiment with four maximum n -gram sizes:

1. **5-gram HHM**: an 11 word window (5 words to the left and right of the target) where the model learns all 2- to 5-grams within that window,
2. **11-gram HHM**: an 11 word window where the model learns all 2- to 11-grams within that window,
3. **21-gram HHM**: a 21 word window where the model learns all 2- to 21-grams within that window, and
4. **Sentence HHM**: a sentence-length window, where the model learns all bigrams to sentence length n -grams within that window.

For all models, the window cannot cross sentence boundaries (e.g., in a five-word sentence, the 21-gram HHM uses a five-word window). Note that for the 5-gram HHM, the maximum n -gram size (5) is distinct from the window size (11), whereas for the three other models the window size is also the maximum n -gram size. We consider large window sizes to account for human sensitivity to long-range dependencies in language, though given that humans can, in

473 principle, be sensitive to arbitrarily long-range dependencies, we consider the
474 fixed context window to be an approximation (see §5.2 for discussion).

475 We use the four HHMs for the following experiments.

476 3.3. *Experiment 1: Part of Speech*

477 If higher-order associations are useful for knowing how a word can be
478 appropriately used in a grammatical sentence, we should expect to see that
479 higher orders of associations enhance the sensitivity of the model to mea-
480 sures of how words are used. In this section, we explore correlations between
481 HHM’s representations and part-of-speech (noun, verb, adverb, adjective,
482 etc.). In the next section, we examine the correlation between HHM’s rep-
483 resentations and the syntactic types proposed by Combinatory Categorical
484 Grammar (CCG; Steedman and Baldridge, 2011).

485 Using WordNet (Princeton University, 2010) and the Moby Part-Of-
486 Speech List (Ward, 1996), we assign a set of part-of-speech tags to each
487 word in the 39 076 word vocabulary. We use similarity between words that
488 are the same part-of-speech as a proxy measure for knowledge that those
489 words can be used in similar ways.

490 Properly speaking, part-of-speech is a theory of language, rather than a
491 behavioral phenomenon, and as such, a cognitive model of language use need
492 not account for part-of-speech as long as it can account for how humans
493 produce and comprehend sentences. Nevertheless, looking at the relation-
494 ship between the representations of HHM and part-of-speech categories can
495 illustrate the effect of the higher levels of the model.

496 Here we analyze the 11-gram HHM, as it is the model with the highest
497 correlation to CCG types in Experiment 2. However, we inspected other
498 window sizes for this analysis and did not observe substantive differences.
499 To examine the effect of higher-order associations, we compare Levels 1 and
500 2 (i.e., second- vs fourth-order associations), Levels 2 and 3 (i.e., fourth vs
501 sixth), and Levels 3 and 4 (i.e., sixth vs eighth).

502 To provide clear examples of higher-order associations and their relation-
503 ships to part-of-speech, we limit our initial analysis to words with at least
504 1000 occurrences in the corpus, as these words have the most robust vec-
505 tor representations. While the part-of-speech of some words (e.g., *manager*,
506 a noun) may be easy to learn from only a few examples, words with more
507 flexible part-of-speech (e.g., *course*, which can be used as a noun, verb, or
508 adjective) may require more examples to learn all the ways in which the word

509 can used, particularly if one of the uses is obscure (e.g., *entire*, *n.*, an uncastrated horse). We also limit our initial analysis to the 500 unique word pairs
510 that increase or decrease in similarity the most between levels. By unique,
511 we mean that we select word pairs where neither word in the pair is present
512 in any of the other 500 word pairs to ensure statistical independence.
513

514 By limiting our analysis in these ways, we focus on unambiguous examples
515 of relationships between words that are affected by fourth-order associations.
516 However, by limiting our analysis, we limit the scope of our conclusions in
517 this analysis to high-frequency words and strong associations. As such, we
518 also conduct more general analyses in this and later sections.

519 To illustrate the nature of higher-order associations, the word pairs that
520 changed the most in similarity between pairs of levels are shown in Table 2.
521 The word pairs that increase the most from Level 1 to 2 can be understood
522 as the most pure examples of words with fourth-order associations but no
523 second-order associations. For example, *focusing* and *derived* have a cosine
524 of -0.10 at Level 1, indicating no second-order association, but have a cosine
525 of 0.86 at Level 2, indicating a strong fourth-order association. Likewise, the
526 word pairs that increase the most from Level 2 to 3 can be understood as
527 the most pure examples of sixth-order associations, and from Level 3 to 4,
528 eighth-order associations.

529 We can see in Table 2, that the four word pairs that increase the most in
530 similarity from Level 1 to 2 are unrelated in meaning, which suggests that
531 second-order associations are sufficient for semantics. However, the top four
532 word pairs that increased the most in similarity from Level 1 to 2 each have
533 exactly matching part-of-speech. While the words *focusing* and *derived* are
534 unrelated in meaning, they are both typically verbs that can also be used as
535 adjectives (e.g., a *focusing lens* or a *derived equation*). Likewise, *focused* and
536 *emerging* can both be used as either an adjective or a verb.

537 By contrast, from Level 3 to 4, the word pair that increases the most in
538 similarity is *across* and *druid*, which has neither meaning nor part-of-speech
539 in common. The word pairs that increase and decrease the most from Level 3
540 to 4 suggest that Level 4 may not provide useful linguistic information.

541 From Level 1 to Level 2, the three word pairs that decrease the most in
542 similarity have partially matching part-of-speech: *clerk* and *local* can both
543 be used as nouns (e.g., *local* in the sense of a *local union branch*), as can
544 *manager* and *main* (e.g., *main* as in a *water main*), and *operator* and *entire*
545 (i.e., *entire* as in an *uncastrated horse*). However, the use of *local*, *main*,
546 and *entire* as nouns is highly infrequent, whereas each is commonly used as

Table 2: The four word pairs that increased or decreased the most in similarity between each level, with each word’s parts of speech (POS) and each word pairs’ change in cosine similarity between levels ($\Delta \cos$). Matching part-of-speech in **bold**.

levels	word 1	word 2	POS 1	POS 2	$\Delta \cos$
1 to 2	focusing	derived	<i>v.</i> , <i>adj.</i>	<i>v.</i> , <i>adj.</i>	+0.95
	searching	associated	<i>v.</i> , <i>adj.</i>	<i>v.</i> , <i>adj.</i>	+0.93
	focused	emerging	<i>v.</i> , <i>adj.</i>	<i>v.</i> , <i>adj.</i>	+0.92
	perched	emerged	<i>v.</i>	<i>v.</i>	+0.92
	clerk	local	<i>n.</i> , <i>v.</i>	<i>n.</i> , <i>adj.</i>	-0.37
	manager	main	<i>n.</i>	<i>n.</i> , <i>adj.</i>	-0.37
	operator	entire	<i>n.</i>	<i>n.</i> , <i>adj.</i>	-0.37
	truth	outer	<i>n.</i>	<i>adj.</i>	-0.37
2 to 3	beings	accord	<i>n. plural</i>	<i>n.</i> , <i>v.</i>	+0.55
	course	cent	<i>n.</i> , <i>v.</i> , <i>adv.</i>	<i>n.</i>	+0.50
	lone	amounts	<i>adj.</i>	<i>n. plural</i> , <i>v.</i>	+0.50
	prime	bye	<i>n.</i> , <i>v.</i> , <i>adj.</i>	<i>exclam.</i> , <i>n.</i>	+0.48
	eh	velvet	<i>exclam.</i>	<i>n.</i> , <i>adj.</i>	-0.14
	huh	silk	<i>exclam.</i>	<i>n.</i>	-0.12
	creaked	hemisphere	<i>v.</i>	<i>n.</i>	-0.11
	erupted	regions	<i>v.</i>	<i>n. plural</i>	-0.11
3 to 4	across	druid	<i>prep.</i> , <i>adv.</i>	<i>n.</i>	+0.37
	course	ought	<i>n.</i> , <i>v.</i> , <i>adv.</i>	<i>n.</i> , <i>v.</i> , <i>adv.</i>	+0.37
	been	Russians	<i>v.</i>	<i>n. plural</i>	+0.36
	must	fraction	<i>n.</i> , <i>v.</i>	<i>n.</i> , <i>v.</i>	+0.36
	huh	which	<i>exclam.</i>	<i>det.</i> , <i>pron.</i>	-0.05
	eh	however	<i>exclam.</i>	<i>adv.</i>	-0.04
	distinction	nineteenth	<i>n.</i>	<i>n.</i> , <i>adj.</i>	-0.03
	but	furthermore	<i>adv.</i>	<i>adv.</i>	-0.03

547 adjectives. As such, these three word pairs are better understood as examples
548 of mismatching part-of-speech (nouns vs. adjectives). Because a partial part-
549 of-speech match is not indicative of the relative frequency of the multiple uses
550 of the word, it is difficult to interpret whether a partial match is more like a
551 match or a mismatch. Thus, we focus our analysis on exact matches.

552 For the 500 word pairs that increased or decreased the most in similar-

ity between each level, Figure 3 shows how many are exact part-of-speech matches, partial matches, or have mismatching part of speech. In total, 13% of all words pairs in the lexicon are exact part-of-speech matches. Among the 500 unique word pairs that increased the most from Level 1 to Level 2, there are significantly more (18%) exact matches than would be expected in a random sample of word pairs ($p < 0.01$). Of the 500 unique word pairs that decreased in similarity the most from Level 1 to 2, 9% are exact matches (e.g., both *great* and *stranger* can be used as an adjective and a noun), which is significantly fewer than expected in a random sample ($p < 0.01$).

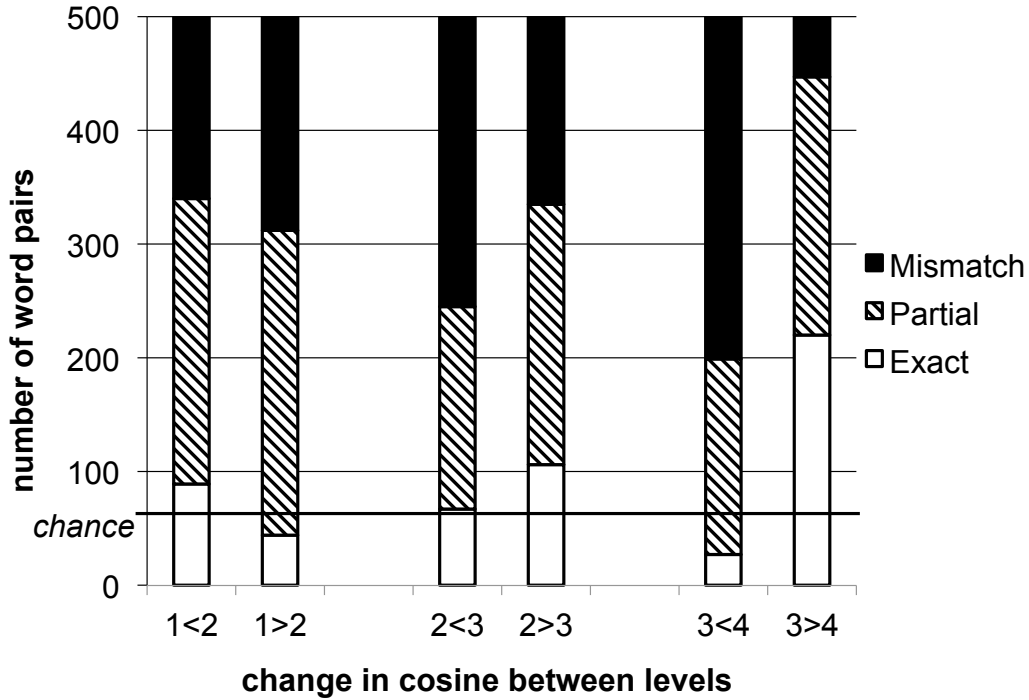


Figure 3: 500 unique word pairs that increased/decreased the most in similarity at each level, categorized by part-of-speech match.

However, from Level 2 to 3 and from Level 3 to 4, significantly more exact matches than expected in a random sample are among the top 500 word pairs that decrease the most ($p < 0.0001$). From Level 3 to 4, significantly fewer exact matches are among the 500 word pairs that increase the most ($p < 0.0001$).

567 The reversal suggests that fourth-order associations are sufficient to dis-
568 cover most exact part-of-speech matches. Indeed, from Level 2 to 3, among
569 the 500 word-pairs that decrease the most, the exact matches have a mean
570 decrease in similarity of 0.00, with a mean cosine of 0.90 between the words
571 at both Levels 2 and 3. Likewise, from Level 3 to 4, the mean decrease
572 in similarity is 0.00 for exact matches, with a mean cosine of 0.98 at both
573 Levels 3 and 4. The exact matches are already highly similar by Level 2
574 and remain highly similar at Levels 3 and 4, and as such their similarity is
575 increased little by sensitivity to sixth- and eighth-order associations.

576 Our analysis thus far has focused on a highly select sample of the corpus:
577 words that occur at least 1000 times and word pairs whose similarity changes
578 dramatically between levels. For the purposes of a more general analysis, we
579 test the ability of HHM to classify the parts-of-speech of all words that occur
580 in the corpus at least five times. Among the 37 543 words that occur at least
581 five times, there are 104 unique sets of part-of-speech tags. We construct a
582 prototype for each part-of-speech tag set as a sum of the vectors for each
583 word that has the exact same tag set. We then classify each word in the
584 lexicon according to the closest prototype, as measured by cosine similarity.

585 As shown in Figure 4a, at Level 1, 20% of words are closest to the proto-
586 type that matches the word’s parts-of-speech. Classification accuracy mod-
587 estly improves at Levels 2 (22%) and 3 (23%) before declining at Level 4
588 (19%). Accuracy is not high, however, as there are 104 part-of-speech pro-
589 totypes, chance classification accuracy is at 1% correct.

590 We re-run the classification only using each word’s most frequent part-of-
591 speech tag in the corpus. We identify the dominant tag using the Stanford
592 Log-Linear Part-of-Speech Tagger (Toutanova et al., 2003) from the Stanford
593 CoreNLP package (Manning et al., 2014)⁴. Again, we exclude words that
594 occur less than five times in the corpus, as well as words with unique part-of-
595 speech tags (e.g., the word “to” is the only word assigned the tag “TO”) for a
596 total of 37 539 words and 29 part-of-speech tags. We compute a prototype for
597 each of the 29 tags and assign words to the closest tag. Chance classification
598 accuracy is 3%.

599 As shown in Figure 4b, at Level 1, 53% of word are classified correctly.
600 Accuracy increases to 62% at Level 2, plateaus at Level 3 (61%) and decreases
601 at Level 4 (58%). Misclassifying nouns and adjectives as each other is the

⁴<https://stanfordnlp.github.io/CoreNLP/>

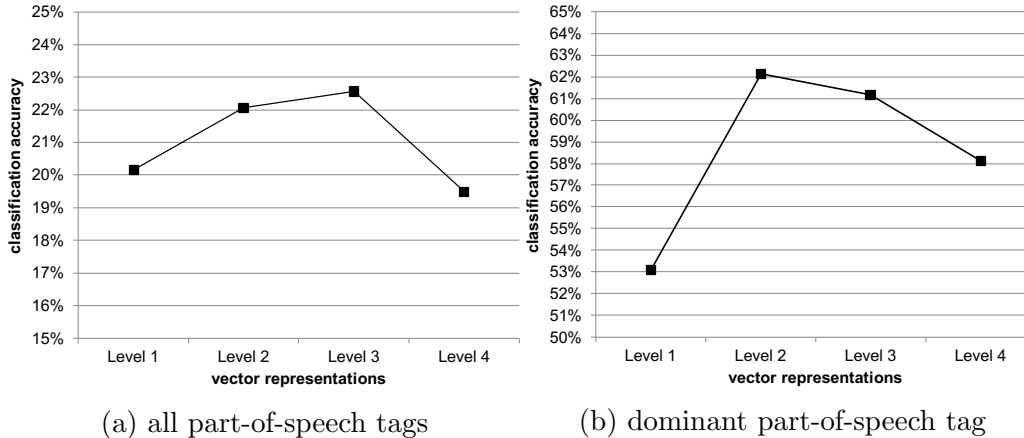


Figure 4: Classification accuracy of words by closest part-of-speech prototype.

largest single source of classification errors at Level 1. At Level 1, 14% of all classifications are errors from confusing nouns and adjectives, compared to only 5% of all classifications at Level 2. The gain in classification accuracy from Level 1 to 2 is mostly due to correctly distinguishing adjectives and nouns. Conversely, confusing singular and plural nouns is a source of error across all levels (7% of all classifications at Level 1 vs. 10% at Level 4), likely due to HHM’s insensitivity to case marking (see §5.4 for discussion).

In summary, strong fourth-order associations (Level 2) strengthen similarities between words with matching part of speech and weaken similarities between words with mismatching part of speech. However, sixth- and eighth-order associations (Levels 3 and 4) do little to further increase similarity between words with the same part-of-speech, and eight-order (Level 4) associations may even obfuscate part-of-speech information.

3.4. Experiment 2: Combinatory Categorical Grammar

Part-of-speech categories (nouns, verbs, adjectives, adverbs) provide a coarse-grained analysis of how words are used in English. Combinatory Categorical Grammar (CCG; Steedman and Baldridge, 2011) is a theory of grammar that provides a more fine-grained analysis of how words are used.

In CCG, sentences are constructed by combining words using a small number of very simple rules. The complexity of language arises not from the complexity of the rules, but from the complexity of the words in the language. In CCG, there are hundreds of types of words, and the type of the

word determines how it can be combined with other words.

The high dimensional space of HHM provides a rich representation of how a word is used in language. As such, correlation between HHM space and CCG type may be more informative than correlation between HHM space and part-of-speech categories.

To classify the words in HHM by CCG type, we use the *Switchboard* Corpus (Godfrey et al., 1992). The Switchboard Corpus is a collection of 2500 telephone conversations. The syntactic structure of the corpus has been annotated using CCG (Reitter et al., 2006). There are 10 256 unique words in the corpus. Of those words, we use the 8768 words that are also in the Novels Corpus. Just as a word can be both an adjective and a verb, a word can have multiple CCG types. To represent the CCG type profile of a word, we represent each word in the Switchboard Corpus by a vector of 357 dimensions, one dimension for each CCG type in the corpus, where the value in each dimension is a count of the number of times that word appears as the given CCG type in the corpus.

The CCG type vectors define similarity relationships between the set of 8768 words. We compute a 8768 x 8768 similarity matrix by taking the cosine of each pair of vectors. To compare relationships in CCG space to relationships in HHM space, we also compute a 8768 x 8768 similarity matrix for each level of HHM. To measure the correlation between CCG space and HHM spaces, we use Spearman’s rank correlation coefficient, which is a non-parametric measure of monotonic (linear or non-linear) relationships in data.

We compute the Spearman’s correlation between the CCG cosine matrix and the cosine matrix for each level of HHM. Figure 5 shows the correlation for each Level of HHM and each maximum n -gram size. The 11-gram HHM achieved the highest correlation with CCG types across all levels, peaking at Level 3 with a correlation of 0.382.

A correlation of 0.382 is not especially high, but it is worth noting that a low correlation does not indicate that HHM is wrong or that CCG is wrong. HHM’s representations contain semantic information that CCG types do not contain. Likewise, CCG types may contain some particulars of syntax that it may be difficult for HHM to learn from a corpus using a sliding context window (see §5.2 for a discussion of using a memory model instead of a fixed window). Other differences may arise simply from how words are used in the Switchboard Corpus versus the Novels Corpus.

HHM’s correlation to CCG is worse when the model includes up to 21-grams or full-sentence-grams, or when restricting the model to 5-grams.

662 Larger n -grams are not always better: as larger n -grams are more unique,
 663 they may be less useful for making inferences about new sentences.

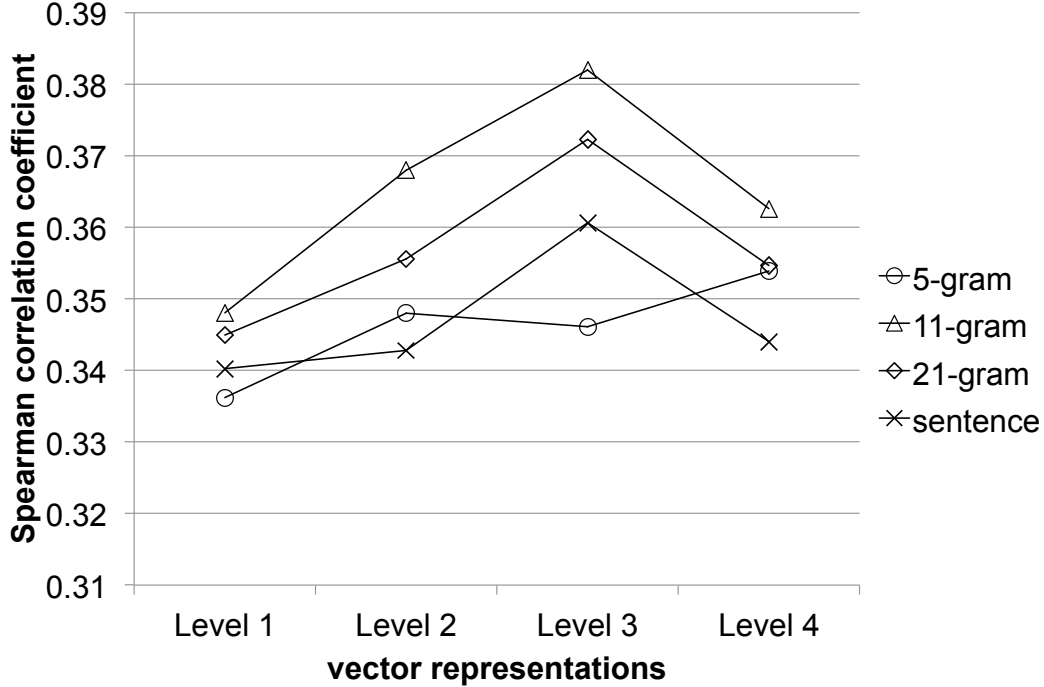


Figure 5: Spearman’s rank correlation coefficient between HHM vectors and CCG types.

664 We see the same pattern across the four levels for the 11-gram, 21-gram,
 665 and full-sentence HHM, with increasing correlation until Level 3, and then a
 666 decrease at Level 4. Though correlation to CCG types is lowest at Level 1
 667 for all models, the increase in correlation is modest, indicating that Level 1
 668 can account for much of the information captured by CCG types.

669 The 5-gram model does not replicate this pattern of correlation, which
 670 we attribute to differences in how the 5-gram model is constructed. Whereas
 671 the other models compute all n -grams within the window, the 5-gram model
 672 computes only 2 to 5-grams within an 11-word window. The dissociation
 673 between window size (11) and maximum n -gram (5), appears to produce a
 674 different behavioural profile than when window size and maximum n -gram
 675 size are scaled together, though we do still see a general upward trend in
 676 correlation at higher levels. The 5-gram model is essentially just the 11-
 677 gram model with the 6- to 11-grams removed, as both models look forward

678 and backward 5 words from the target word. However, the 5-gram model’s
679 correlation to CCG types is lower than the 11-gram model at all levels, which
680 suggests that the 6- to 11-grams are, in fact, providing useful information
681 about syntax, independently from the size of the context window.

682 In summary, higher-order associations (up to Level 3, i.e., sixth-order
683 associations) improve the ability of the model to capture syntactic type re-
684 lationships and that large n -grams, in the range from 6-grams up to at least
685 11-grams, provide useful information about the syntactic type of words.

686 3.5. Experiment 3: Word Ordering Task

687 The real test of syntactic knowledge is the ability to form grammatical
688 sentences. Do higher-order associations provide additional useful information
689 about how to sequence words into a grammatical sentence? When given an
690 unordered set of words that can be arranged into a sentence, are higher levels
691 of HHM better able to find the grammatical ordering?

692 We replicate a task from Johns et al. (2016a). In this task, a model is
693 given an unordered set of n words taken from an n -word sentence. The model
694 must discern which of the $n!$ possible word orderings is the original ordering.

695 HHM is not, by itself, able to perform the word ordering task, because
696 HHM does not operate on sentences. However, HHM’s representations con-
697 tain word-level information that can be leveraged to perform the task when
698 situated within a sentence-level model. We use a simplified version of the ex-
699 emplar model used by Johns et al. (2016a). The exemplar model is provided
700 with an exemplar set consisting of 125 000 seven-word sentences randomly
701 sampled from the Novels Corpus. Sentences in the exemplar set have no
702 words with frequency less than 300. All test set sentences and permutations
703 thereof are excluded from the exemplar set.

704 We embed the word representations generated by each level of HHM in
705 the exemplar model. Each sentence in the exemplar set is represented as
706 a pair of vectors in the exemplar model. One vector is an unordered set
707 of words constructed as a sum of HHM’s memory vectors representing each
708 word in the sentence. The second vector is the sum of all ordered sequences of
709 words in the sentence, from individual words up to 7-grams. Each sequence
710 is constructed as a convolution of HHM’s memory vectors for each word in
711 the sequence. Before use, all HHM vectors are normalized to a Euclidean
712 length of one and permuted, as shown in Equation 2.

713 Test items are a set of 200 seven-word sentences as used by Johns et al.
714 (2016a). Test items have simple syntactic construction and consist of words

715 that occur at least 300 times in the corpus. Test items are presented to the
 716 exemplar model as an unordered set of words.

717 The exemplar model first selects the exemplar sentence most similar to the
 718 test item, as measured by the cosine between the vectors for the unordered
 719 sets. Then, of the $7!$ possible orderings of the words in the test item, the
 720 model selects the ordering most similar to that of the selected exemplar
 721 sentence, as measured by the cosine between the vectors representing the
 722 ordered sequences of words. The ordering produced by the model is judged
 723 to be correct if it matches the original ordering of the words in the test item.

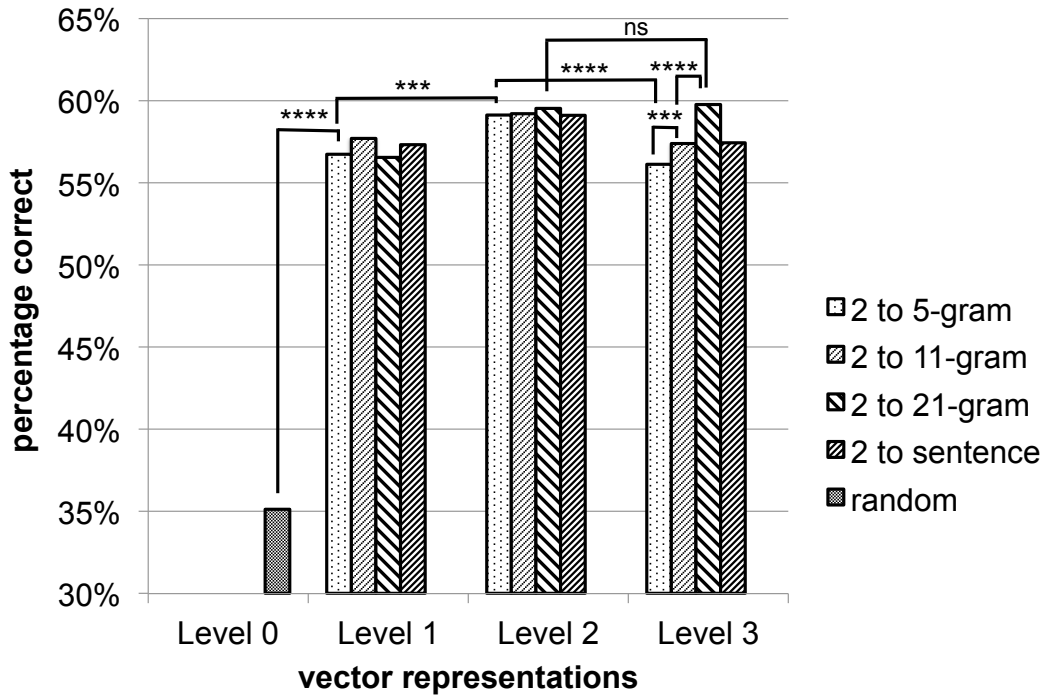


Figure 6: Test sentences correctly ordered by model as a function of vectors used to represent words.

724 We test all four versions of HHM from Level 1 to Level 3. To ensure that
 725 results are not contingent on a particular sample of 125 000 exemplar sen-
 726 tences, results are averaged across 50 random samples. Mean percent correct
 727 across the 50 samples is shown in Figure 6. To test for statistical significance
 728 across the seven conditions, we use a repeated-measures permutation test, a
 729 non-parametric measure (Mewhort et al., 2010, 2009).

730 We also include a “Level 0” as a baseline for performance. Level 0 rep-
731 represents individual words as randomly generated vectors and the sentence
732 vectors are constructed from those vectors. In effect, at Level 0, the model
733 selects the exemplar sentence with the most words in common with the test
734 item and applies the word ordering of the selected exemplar to the test item.
735 Level 0 provides a baseline where the model is sensitive to neither seman-
736 tic similarity nor higher-order associations but is sensitive to word overlap
737 between the test item and exemplars. Level 0 gets a mean of 35.1% correct.

738 Level 1 outperforms Level 0 across all window sizes ($p < 0.0001$) with a
739 mean of 57.1% correct. Level 1 selects the exemplar sentence that has the
740 most semantic similarity to a given test item.

741 Level 2 outperforms Level 1 across all window sizes ($p < 0.001$) with
742 a mean of 59.2% correct, demonstrating that fourth-order associations con-
743 tribute to the task of ordering words into grammatical sentences.

744 At Level 3, performance declines for all models $p < 0.0001$ except the
745 21-gram HHM, for which performance does not change significantly from
746 Level 2 to 3 ($p > 0.05$). Here we see a significant effect of window size. The
747 21-gram HHM outperforms all other Level 3 models ($p < 0.0001$) and the
748 5-gram HHM performs worse than all other Level 3 models ($p < 0.0001$).

749 Inspecting by hand the errors made in a single run of Level 0 and each
750 level of the 21-gram HHM, we find that the pattern of errors varies little
751 across levels. The Levenshtein edit distance from a produced error to a cor-
752 rect ordering has a mean of 3 at each level of the model. All levels occasion-
753 ally suggest a grammatical ordering of the words different from the original
754 ordering (e.g., “he opened the door and got up”, an incorrect ordering at
755 Level 3, is grammatical even if “he got up and opened the door” would be a
756 more typical sequence of actions). At Level 0, we found that an additional
757 6.5% of the 200 sentences produced were grammatical but not the original
758 ordering. At Level 1, we found an additional 11.5% to be grammatical, at
759 Level 2, an additional 11.0%, and at Level 3, an additional 7.5%. The re-
760 maining incorrect orderings are ungrammatical, typically due to a misplaced
761 word (e.g., “came a serious look over his face”, at Level 1, misplaces the verb
762 *came*, or “I do not much trust you that”, at Level 2 misplaces *much*).

763 Our results show that for the task of ordering words into grammatical
764 sentences, a model that uses fourth-order associations between words out-
765 performs a model that uses second-order associations. Our results also show
766 that a model that uses second-order associations or higher outperforms a
767 model that only uses word overlap (i.e. Level 0).

768 The results show little benefit to using a window beyond 5-grams, possibly
769 because the task is restricted to constructing 7-gram sentences. However, the
770 5-gram HHM performs the worst at Level 3 and the 21-gram HHM performs
771 the best, which suggests there are two counter-acting processes at work. At
772 higher levels, HHM is increasingly able to make useful inferences about the
773 relationships between large, low frequency n -grams, while simultaneously
774 losing the ability to make fine discriminations between small, high frequency
775 n -grams. We hypothesize that the decline in task performance from Levels
776 2 to 3 is due to all HHMs losing the ability to make fine discriminations
777 for small n -grams. Performance of HHM representations that contain larger
778 n -grams is less affected as those models are simultaneously gaining an ability
779 to better use those large n -grams.

780 To test this hypothesis, we break down HHM into its constituent n -gram
781 components. While the HHMs previously discussed learned 2-grams up to
782 n -grams for some n , here we train each HHM on one and only one size of
783 n -gram. For Level 0, we use random vectors. For Level 1, we use Level 0’s
784 random vectors to construct a 2-gram only HHM, a 3-gram only HHM, etc.,
785 up to a 7-gram only HHM. For Level 2, we construct the HHMs out of Level 1
786 of the 2- to 21-gram HHM. For Level 3, we construct the HHMs out of Level 2
787 of the 2- to 21-gram model. We use the 2- to 21-gram HHM as it is the model
788 with the most robust performance across all levels on this task.

789 Figure 7 shows the percentage of test sentences ordered correctly across
790 different sizes of n -gram and levels of HHM. Results are averaged across 10
791 sets of 125 000 exemplar sentences. Higher levels of the model are better able
792 to use larger n -grams. Level 1 of HHM is best able to use 2-grams and 3-
793 grams. Conversely, at Level 3 of HHM, the model is able to make use of large
794 n -grams, but performance declines for smaller n -grams. Task performance
795 at Level 2 of HHM peaks for 3- to 6-grams, and declines for 2- and 7-grams.
796 Level 0 is included for a baseline performance of 35.1% correct.

797 Figure 7 illustrates that at higher levels, HHM progressively loses the
798 ability to make fine distinctions between small n -grams as the representations
799 for the words that compose the n -grams become increasingly similar. For
800 example, “she grinned” and “he smiled” may be represented by identical or
801 nearly identical bigrams at higher levels. However, higher levels begin to be
802 able to make use of large n -grams. At lower levels, large n -grams are unique,
803 and thus do not provide useful information about the relationships between
804 words. At higher levels, large n -grams are similar to other large n -grams.
805 For example, while the 7-gram “you are as gregarious as a locust” may occur

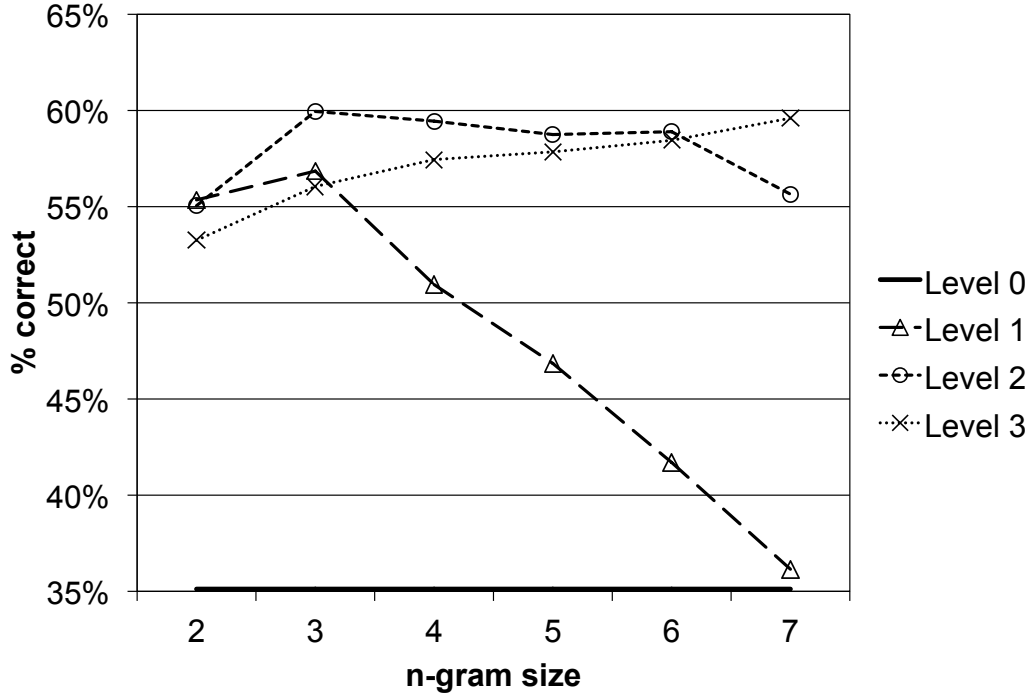


Figure 7: Test sentences correctly ordered as a function of n -gram size and HHM level.

only once in a corpus, at higher levels of HHM, this 7-gram comes to resemble other 7-grams, such as “he was as strong as an ox”.

Correctness is a noisy metric of model skill as it is binary. We can get a more precise measure of model skill by using the cosine scores assigned to each of the $7!$ alternative orderings. To measure the degree of confidence with which the model endorses a given ordering as grammatical, we use the deviation of the grammatical ordering’s cosine from the cosines of the other orderings. The deviation is a graded measure, sensitive to how close the model is to wrong when it’s correct and how close to correct the model is when it’s wrong, giving us a better picture of the model’s decisions. We normalize the deviation by the standard deviation to control for differences in the spread of cosine values at different levels.

As shown in Figure 8, the deviations yield the same pattern of results as Figure 7. The ability of the Level 1 models to discriminate between the correct answer and alternatives is highest for 2-grams and 3-grams and declines

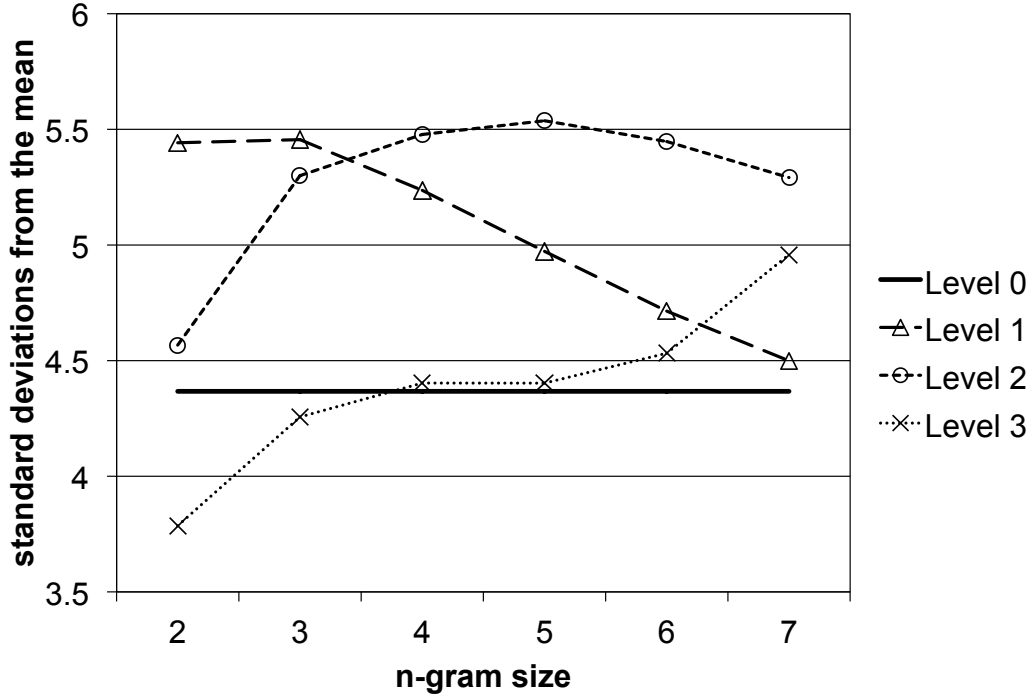


Figure 8: Deviation of correct word ordering from alternatives as a function of n -gram size and HHM level.

for larger n -grams. At Level 3, we observe the opposite: the deviation of the correct answer is highest for 7-grams and declines for smaller n -grams. At Level 2, deviation peaks at 5-grams, declining for smaller or larger n -grams.

The results in Figures 7 and 8 demonstrate that the higher levels of HHM allow for better use of large n -gram information, at the cost of a declining ability to make discriminations between small n -grams. Specifically, Level 1 (i.e., second-order associations) representations make the best use of 2-grams and 3-grams, Level 2 (i.e., fourth-order associations) makes the best use of 4-grams to 6-grams, and Level 3 makes the best use of 7-grams.

Note that none of the combined n -gram models in Figure 6 outperform the 3-gram only Level 2 model or the 7-gram only Level 3 model, which suggests that simpler HHMs may be sufficient for the word ordering task.

Accuracy on the task can be increased by using more data. By increasing the size of the exemplar set from 125 000 sentences to 500 000, accuracy for the 21-gram Level 2 HHM improves from 60% correct to 68% correct.

836 Performance can also be improved by using a more complex model. Rather
837 than selecting the best exemplar, Johns et al. (2016a) use a weighted sum of
838 all exemplars to make word-ordering decisions. Each exemplar is weighted
839 by the cosine similarity between the exemplar and the unordered set of words
840 in the test item, raised to a power (a fitting parameter).

841 Johns et al. (2016a) report a best accuracy of 60% with 500 000 exemplars,
842 random vectors (Level 0), and an exponent of 9. We find a best accuracy of
843 76% correct with 500 000 exemplars, the 21-gram Level 2 HHM vectors, and
844 an exponent of 450. However, our aim is not to optimize accuracy on the
845 word-ordering task, but to illustrate the role of higher-order associations in
846 constructing grammatical sentences.

847 3.6. Experiment 4: *Colorless green ideas sleep furiously*

848 Chomsky (1956) gives “Colorless green ideas sleep furiously” as an exam-
849 ple of a sentence that is grammatically correct but meaningless. By contrast,
850 Chomsky notes that “Furiously sleep ideas green colorless” is ungrammati-
851 cal. Chomsky uses this example as an argument against statistical models of
852 speech. Unless the sentence “Colorless green ideas sleep furiously” is part of
853 the statistical model’s training corpus, a statistical model would neither be
854 able to generate the sentence nor determine that it is grammatical.

855 Pereira (2000) demonstrates that a statistical model can, in fact, discrim-
856 inate between “Colorless green ideas sleep furiously” and the ungrammati-
857 cal “Furiously sleep ideas green colorless”. Pereira uses an *aggregate bigram*
858 *model* that estimates the probability of each bigram in “Colorless green ideas
859 sleep furiously” by using second-order associations to known bigrams and an
860 expectation-maximization algorithm (Dempster et al., 1977). Pereira’s ag-
861 gregate bigram model finds that “Colorless green ideas sleep furiously” is
862 about 20 000 times more likely than “Furiously sleep ideas green colorless”.

863 HHM is also a statistical model that can be understood as estimating the
864 probability of unseen n -grams through the use of higher-order associations.
865 Can the higher levels of HHM discern that “Colorless green ideas sleep fu-
866 riously” is a grammatical sentence? Given the unordered set of five words
867 *colorless*, *furiously*, *green*, *ideas*, and *sleep*, there are $5! = 120$ possible order-
868 ings of those words. Does HHM demonstrate a better than chance preference
869 for Chomsky’s grammatical but meaningless ordering of the words over “Fu-
870 riously sleep ideas green colorless” or the 118 other orderings? If HHM is
871 purely semantic and “Colorless green ideas sleep furiously” is a purely syn-
872 tactic sentence, performance should be no better than chance at this task.

873 We use the same exemplar model as in the previous section. To construct
 874 the exemplar model’s vectors, we use the 21-gram HHM. The exemplar model
 875 is provided a set of five-word sentences and picks the sentence most similar to
 876 the unordered set of words *colorless*, *furiously*, *green*, *ideas*, and *sleep*. The
 877 selected sentence’s structure is then used to score the 120 possible orderings.

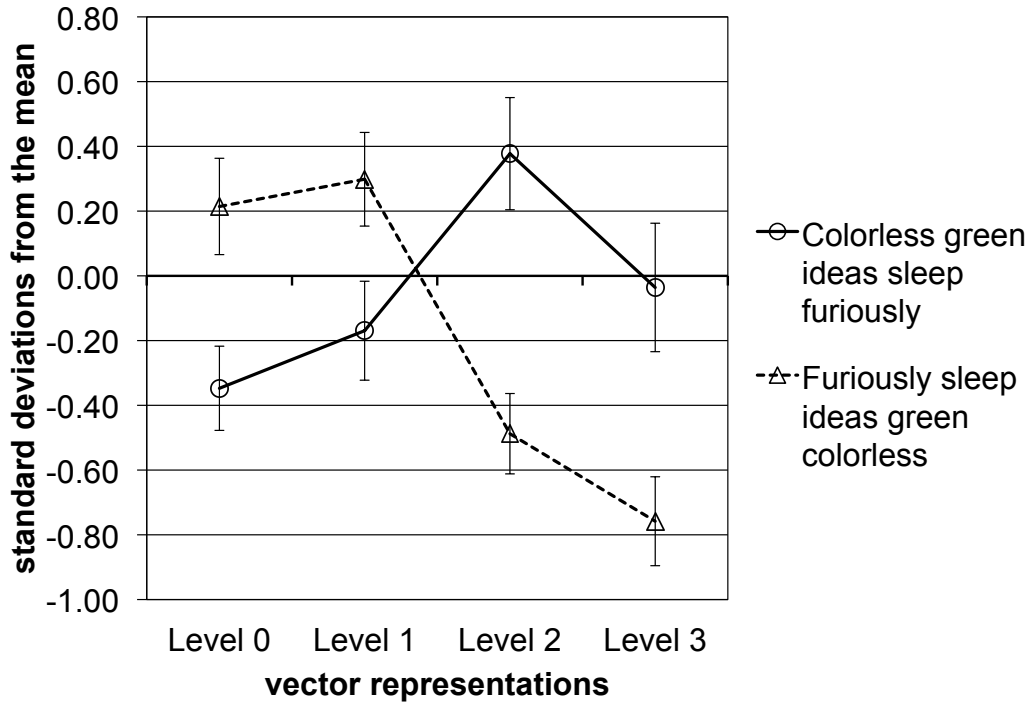


Figure 9: Deviation from the mean cosine score as a function of HHM level.

878 Mean deviation of both “Colorless green ideas sleep furiously” and “Furi-
 879 ously sleep ideas green colorless” at each level of HHM is shown in Figure 9.
 880 Results are averaged across 50 different random sets of 125 000 sentences.
 881 Error bars indicate standard error. Word orderings above zero are judged to
 882 be more grammatical than the mean of the 120 possible sentence orderings
 883 and orderings below zero are judged to be less grammatical than the mean.

884 Is “Colorless green ideas sleep furiously” more grammatical, according to
 885 the exemplar model, than “Furiously sleep ideas green colorless”? To test
 886 for statistical significance, we use a repeated-measures permutation test. At
 887 Level 0, “Furiously sleep ideas green colorless” is more grammatical ($p <$

0.05) whereas at Levels 2 and 3, “Colorless green ideas sleep furiously” is more grammatical ($p < 0.05$). Given that the two sentences are exact reverse orderings of each other, it is not surprising that the model’s confidence in each sentence is roughly the inverse of the other’s.

Thus, selecting exemplar sentences with words in common with the test set (e.g., *green*, *furiously*, etc.), as the model does at Level 0, is not enough to make the correct grammatical distinction. Selecting sentences with similar meanings (e.g., *red*, *angrily*, etc.), as Level 1 does, is likewise insufficient. Higher-order associations at Levels 2 and 3, seem to be necessary to determine that “Colorless green ideas sleep furiously” is the more grammatical alternative of the pair.

Identifying “Colorless green ideas sleep furiously” as the *most* grammatical ordering of the 120 possible orderings is a more difficult problem. Only at Level 2 does the model judge “Colorless green ideas sleep furiously” to be more likely than average ($p < 0.05$), selecting it as the most likely ordering 7 times out of 50. The rate at which Level 2 selects “Colorless green ideas sleep furiously” as the preferred alternative might be improved by either increasing the number of exemplars over the current 125 000 or by using a more sophisticated model (e.g., Johns et al., 2016a; Gulordava et al., 2018).

The results suggest that “Colorless green ideas sleep furiously” cannot be judged as grammatical by analogy to sentences with either the same words or words with similar meanings. However, sensitivity to fourth-order associations causes representations for words with similar syntactic type to look increasingly alike, such that “Colorless green ideas sleep furiously”, or *adjective adjective noun verb adverb*, begins to look like an English sentence.

4. Other Models of Higher-Order Associations

While we have based HHM on BEAGLE, it is possible to use other models to detect higher-order associations in language.

The Associative Smoothing Network (Roberts and Chater, 2008), for example, is a spreading activation model that uses third-order associations to make sentence acceptability judgments. The network has no inherent limitation to how far activation can spread, and so can be easily applied recursively to detect fourth-, fifth-, sixth-order associations and higher.

However, for some models, there’s no trivial way to recursively apply the model to incorporate higher-order associations. For example, the *word2vec* neural network expects, for each word it takes as input, a vector that uses

924 *one-hot encoding*⁵. Conversely, the semantic vectors *word2vec* generates have
925 d dimensions, where d is much smaller than the size of the lexicon, and each
926 dimension is real valued and individually meaningless.

927 Because semantic vectors and one-hot vectors have such different proper-
928 ties, the semantic vectors cannot be re-used as input to *word2vec* to recur-
929 sively detect higher-order associations. While it’s almost certainly possible
930 to design a neural network model of distributional semantics that can be
931 recursively applied in much the same manner as HHM, *word2vec* cannot be
932 used to do so as standardly implemented.

933 BEAGLE and HHM also have a unique property that may make repli-
934 cating our results with other models difficult. Other models of distributional
935 semantics only learn relationships between pairs of words, whereas BEAGLE
936 and HHM learn a relationship between a word and sequences of words.

937 Models limited to knowing the relationships between pairs of words can
938 certainly benefit from third- or fourth-order associations. The Associative
939 Smoothing Network, for example, is strictly a bigram model, but third-order
940 associations allow the model to make judgments about the acceptability of
941 novel word pairs (Roberts and Chater, 2008). However, in the word ordering
942 task, we find that improvements in performance at higher orders of associ-
943 ation largely result from improving the ability of HHM to make use of the
944 information in larger n -grams, $n \geq 3$ (see Figures 7 and 8).

945 While we are not committed to the specific implementation details of how
946 the Hierarchical Holographic Model learns higher-order associations, HHM
947 has two desirable properties for modelling higher-orders of association:

- 948 1. HHM can be recursively applied an arbitrary number of times to learn
949 arbitrarily high orders of association, and
- 950 2. HHM is able to learn arbitrarily large n -grams with linear time com-
951 plexity and constant space complexity.

952 5. Future Work

953 The Hierarchical Holographic Model (HHM) has a number of limitations,
954 namely, (1) HHM learns batch-style rather than online, (2) HHM’s fixed
955 window is unrealistic, (3) HHM does not combine levels of representation,

⁵In *one-hot encoding*, a vector has one dimension for each word in the lexicon. To represent a word, that word’s dimension is set to 1 and all other dimensions are set to 0.

956 and (4) HHM is applied only to English. HHM also has possible applications
957 beyond what we explore in this paper, such as (5) modelling developmental
958 language acquisition. We discuss each avenue for future research in turn.

959 5.1. *Online Learning*

960 HHM is not an online model of learning. HHM learns each level of rep-
961 resentation sequentially. In keeping with research on the acquisition of first
962 and second order associations in children and adults (McNeill, 1963; Sloutsky
963 et al., 2017), we would expect learning to happen at each order of associa-
964 tion continuously and in parallel. Though we hypothesize that higher-order
965 associations scaffold off lower-order associations, we hold that the scaffold-
966 ing is such that higher-order associations are adjusted as new lower-order
967 associations are learned.

968 HHM may be able to learn all levels in parallel. Doing so will introduce
969 more noise into the higher levels of the model, as early on, the level(s) below
970 will not have stable memory vectors yet, but over a large enough corpus,
971 stable representations should emerge first at Level 1 and then propagate
972 upward to higher levels.

973 5.2. *Window Size*

974 Humans are sensitive to long-range dependencies in language. For ex-
975 ample, in anaphora resolution, readers are able to identify the referent of a
976 pronoun such as *she* even over a large number of intervening words or sen-
977 tences. Readers selectively and strategically maintain pertinent information
978 in memory from much earlier in a sentence, paragraph, or passage (Kintsch
979 and Van Dijk, 1978).

980 We include large window sizes in our simulations as proxies for the ca-
981 pacity of memory to selectively retain long-range information. The sliding
982 context window of HHM is best understood as an inexact proxy for the lin-
983 guistic associations and dependencies available to a reader (or listener) when
984 the target word in a sentence is encoded. However, human memory does not
985 behave like a verbatim list of the last 21 words read (or heard).

986 To build a more detailed model of human sentence processing, we would
987 need to replace the sliding window with a model of the linguistic informa-
988 tion maintained in working memory and stored in long-term memory, as
989 in Kintsch and Van Dijk’s (1978) model of sentence processing. We would
990 also need a model of selective attention to account for what information is

991 retained in long-term memory and maintained in working memory, as in-
 992 formed by the model’s experience of what is likely to be useful for resolving
 993 the syntax and semantics of future utterances. Computational, holographic
 994 approaches to modelling working memory (Franklin and Mewhort, 2015) and
 995 episodic memory (Jamieson and Mewhort, 2011) could potentially be inte-
 996 grated with HHM to provide a more detailed processing model.

997 5.3. *Combining Levels*

998 Gruenenfelder et al. (2016), modeling word association norms, find that
 999 a hybrid model that uses both first- and second-order associations better
 1000 matches human data. We note that on the word ordering task, while, on
 1001 average, Level 2 with any window, or Level 3 with the 21 word window, pro-
 1002 duces the best results, Level 1 often correctly ordered sentences that Levels
 1003 2 or 3 got wrong. Perhaps a model that uses all three levels could outper-
 1004 form a model that uses only one level at a time. A neural network model
 1005 that combines input from varying n -gram sizes and from varying orders of
 1006 association might be able to outperform a neural network that strictly takes
 1007 traditional word embeddings as input. We hypothesize that human memory
 1008 is able to use relations between concepts at varying levels of abstraction as
 1009 needed to meet task demands.

1010 5.4. *Other Languages*

1011 In languages with extensive case marking (e.g., Latin), case markers are
 1012 used to indicate the part-of-speech of a word instead of relying on word
 1013 order, as English does. To learn the case markers, HHM would need to either
 1014 process the corpus parsed into sub-word units (e.g., Cotterell and Schütze,
 1015 2015), splitting off the case marker from the root word, or to use non-random
 1016 environment vectors that represent the orthography of the word, as in Cox
 1017 et al. (2011).

1018 The utility of HHM’s sensitivity to word sequence and higher-order associ-
 1019 ations for modelling case-marked languages is an open question. Case-marked
 1020 languages typically use word sequence to convey non-syntactic information
 1021 (e.g., emphasis or new information), such that while preserving word order
 1022 may not be important for syntax, *per se*, order remains important for con-
 1023 veying meaning. Thus while the type of information captured by HHM’s
 1024 sensitivity to word sequence and abstract associations may differ in case-
 1025 marked languages, we expect that sequence and associations will still play
 1026 an informative role. HHM’s central hypothesis is that human memory has

the capacity for sensitivity to abstract associations, even if those associations are potentially used differently across languages.

5.5. *Language Acquisition*

Children acquire first-order associations earlier in development than second-order associations (Brown and Berko, 1960; Ervin-Tripp, 1970; Nelson, 1977; Sloutsky et al., 2017). Likewise, McNeill (1963) found that when participants are trained on a set of non-words and are tested with a free association task, after 20 trials of training, participants produce only first-order associations between the non-words, but by 60 trials, participants produce both first- and second-order associations.

Sloutsky et al. (2017) propose a neural network model that captures the gradual acquisition of second-order associations contingent on learning first-order associations sensitive to sequential word order, as well as the acquisition of order-independent first-order (syntagmatic) associations. Similarly, the Syntagmatic-Paradigmatic Model (Dennis, 2004, 2005) is a computational model of human memory and language learning that postulates two long-term memory systems: one for sequences and one for (order-independent) relations, which respectively account for knowledge of first-order (syntagmatic) and second-order (paradigmatic) associations

According to Barcel-Coblijn et al. (2012), the point at which a child transitions from speaking in utterances of one or two words to speaking in full sentences is the point at which the child’s knowledge of the relationships between words transitions from a sparsely connected graph to a dense “small world” graph, typical of an adult vocabulary, where all words are several steps from all other words in the graph. We hypothesize that learning longer range connections between words is necessary to construct novel syntactic utterances. We speculate that a model that captures higher-order associations, such as an online variant of HHM that uses both *context* and *order* vectors, and is therefore sensitive to both super-paradigmatic and super-syntagmatic associations, may be able to account for the dynamics of a child’s language learning process.

6. **Conclusions**

We define orders of association and explore the hypothesis that higher-order associations in language capture syntactic relationships between words. We propose a “deep” model of distributional semantics, the Hierarchical

Holographic Model (HHM), sensitive to higher-order associations. We evaluate the correlation between HHM’s representations, part-of-speech, and the lexical syntactic types of Combinatory Categorical Grammar (CCG; Steedman and Baldridge, 2011). We find that strong fourth-order associations are likely to increase similarity between words with the same part-of-speech and decrease similarity between words with mismatching part-of-speech. Fourth- and sixth-order associations increase correlation with CCG type relative to second-order (i.e., paradigmatic) associations.

Fourth-order associations also improve the ability of HHM’s representations to order words into grammatical sentences, including nonsense sentences such as Chomsky (1956)’s “Colorless green ideas sleep furiously”. The usefulness of higher-order associations interacts with the window size of the distributional semantics model, such that larger n -grams require higher orders of association in order to contribute useful information, whereas smaller n -grams are best represented using lower orders of association.

In summary, we find consistent evidence that fourth-order associations (Level 2) provide useful linguistic information of a syntactic character. Conversely, the evidence is mixed for sixth-order (Level 3), and we find no evidence that eighth-order associations (Level 4) are useful for linguistic tasks.

We hypothesize that humans are also sensitive to higher-order associations in non-linguistic domains. Humans have the ability to abstract away from the specifics of an experience (i.e. episodic memories) to infer concepts (i.e., semantic memories) from the patterns that occur across multiple experiences (e.g., Hintzman, 1986). The theoretical claim of HHM is that the pattern inference process is recursive, such that human memory can also infer meta-concepts from patterns across concepts, and that these meta-concepts play an important role in human behaviour, such as language.

7. Acknowledgments

We thank Eilene Tomkins-Flanagan, Kevin D. Shabahang, and the late D. J. K. Mewhort for the use of their BEAGLE code. We are also indebted to D. J. K. Mewhort for comments on the paper and the use of his server, funded by a grant from the Natural Sciences and Engineering Research Council of Canada (NSERC: APA 318). We thank Jeremy R. Cole for use of his CCG vectors and for his comments on the paper. We thank Saranya Venkatraman for her insights regarding the *word2vec* model. This research has been funded by an Ontario Graduate Scholarship and NSERC Post-Doctoral Fellowship

1098 to M. A. Kelly, a National Science Foundation grant (BCS-1734304) to David
1099 Reitter and M. A. Kelly, and a grant from NSERC to Robert L. West.

1100 References

1101 Anderson, J.R., 2009. How can the human mind occur in the physical uni-
1102 verse? Oxford University Press, New York, NY.

1103 Barcel-Coblijn, L., Corominas-Murtra, B., Gomila, A., 2012. Syntactic trees
1104 and small-world networks: Syntactic development as a dynamical process.
1105 Adaptive Behavior 20, 427–442. doi:10.1177/1059712312455439.

1106 Brown, R., Berko, J., 1960. Word association and the acquisition of grammar.
1107 Child Development 31, 1–14.

1108 Burgess, C., Lund, K., 1997. Modelling parsing constraints with high-
1109 dimensional context space. Language and Cognitive Processes 12, 177–210.

1110 Chomsky, N., 1956. Three models for the description of language. IRE
1111 Transactions on Information Theory 2, 113–124. doi:10.1109/TIT.1956.
1112 1056813.

1113 Cotterell, R., Schütze, H., 2015. Morphological word-embeddings, in: Pro-
1114 ceedings of the 2015 Conference of the North American Chapter of the
1115 Association for Computational Linguistics: Human Language Technolo-
1116 gies, Association for Computational Linguistics, Denver, Colorado. pp.
1117 1287–1292. doi:10.3115/v1/N15-1140.

1118 Cox, G.E., Kachergis, G., Recchia, G., Jones, M.N., 2011. Towards a scalable
1119 holographic representation of word form. Behavior Research Methods 43,
1120 602–615. doi:10.3758/s13428-011-0125-5.

1121 Dempster, A.P., Laird, N.M., Rubin, D.B., 1977. Maximum likelihood from
1122 incomplete data via the em algorithm. Journal of the royal statistical
1123 society. Series B (methodological) , 1–38URL: [http://www.jstor.org/
1124 stable/2984875](http://www.jstor.org/stable/2984875).

1125 Dennis, S., 2004. An unsupervised method for the extraction of propositional
1126 information from text. Proceedings of the National Academy of Sciences
1127 101, 5206–5213. doi:10.1073/pnas.0307758101.

- 1128 Dennis, S., 2005. A memory-based theory of verbal cognition. *Cognitive*
1129 *Science* 29, 145–193. doi:10.1207/s15516709cog0000\9.
- 1130 Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2019. BERT: Pre-training
1131 of deep bidirectional transformers for language understanding, in: *Pro-*
1132 *ceedings of the 2019 Conference of the North American Chapter of the*
1133 *Association for Computational Linguistics: Human Language Technolo-*
1134 *gies*, Association for Computational Linguistics, Minneapolis, Minnesota.
1135 pp. 4171–4186. doi:10.18653/v1/N19-1423.
- 1136 Eliasmith, C., 2013. *How to build a brain: A neural architecture for biological*
1137 *cognition*. Oxford University Press, New York, NY.
- 1138 Ervin-Tripp, S.M., 1970. Substitution, context, and association, in: Postman,
1139 L., Keppel, G. (Eds.), *Norms of Word Association*. Academic Press, pp.
1140 383 – 467. doi:10.1016/B978-0-12-563050-4.50012-1.
- 1141 Franklin, D.R.J., Mewhort, D.J.K., 2015. Memory as a hologram: An anal-
1142 ysis of learning and recall. *Canadian Journal of Experimental Psychology*
1143 69, 115–135. doi:10.1037/cep0000035.
- 1144 Gabor, D., 1969. Associative holographic memories. *IBM Journal of Research*
1145 *and Development* 13, 156–159. doi:10.1147/rd.132.0156.
- 1146 Gayler, R.W., 2003. Vector symbolic architectures answer Jackendoff’s chal-
1147 lenges for cognitive neuroscience, in: Slezak, P. (Ed.), *Proceedings of the*
1148 *Joint International Conference on Cognitive Science*. University of New
1149 South Wales, Sydney, Australia, pp. 133–138. URL: <http://cogprints.org/3983/>.
- 1150
- 1151 Godfrey, J.J., Holliman, E.C., McDaniel, J., 1992. Switchboard: Telephone
1152 speech corpus for research and development, in: [Proceedings] ICASSP-
1153 92: 1992 IEEE International Conference on Acoustics, Speech, and Signal
1154 Processing, pp. 517–520 vol.1. doi:10.1109/ICASSP.1992.225858.
- 1155 Grefenstette, G., 1994. Corpus-derived first, second and third-order word
1156 affinities, in: *Proceedings of the Sixth Euralex International Congress*,
1157 *Association for Computational Linguistics*, Amsterdam, The Netherlands.
1158 pp. 279–290.

- 1159 Griffiths, T.L., Steyvers, M., Tenenbaum, J.B., 2007. Topics in semantic
1160 representation. *Psychological Review* 114, 211–244. doi:10.1037/0033-
1161 295X.114.2.211.
- 1162 Gruenenfelder, T.M., Recchia, G., Rubin, T., Jones, M.N., 2016. Graph-
1163 theoretic properties of networks based on word association norms: Im-
1164 plications for models of lexical semantic memory. *Cognitive Science* 40,
1165 1460–1495. doi:10.1111/cogs.12299.
- 1166 Gulordava, K., Bojanowski, P., Grave, E., Linzen, T., Baroni, M., 2018.
1167 Colorless green recurrent networks dream hierarchically, in: *Proceedings*
1168 *of the 2018 Conference of the North American Chapter of the Association*
1169 *for Computational Linguistics: Human Language Technologies, Volume 1*
1170 *(Long Papers)*, Association for Computational Linguistics. pp. 1195–1205.
1171 URL: <http://aclweb.org/anthology/N18-1108>, doi:10.18653/v1/N18-
1172 1108.
- 1173 Hintzman, D.L., 1986. “Schema abstraction” in multiple-trace memory mod-
1174 els. *Psychological Review* 93, 411–428. doi:10.1037/0033-295X.95.4.
1175 528.
- 1176 Jamieson, R.K., Mewhort, D.J.K., 2011. Grammaticality is inferred from
1177 global similarity: A reply to kinder (2010). *The Quarterly Journal of Ex-*
1178 *perimental Psychology* 64, 209–216. doi:10.1080/17470218.2010.537932.
- 1179 Jenkins, J.J., 1964. A Mediational Account of Grammatical Phenomena.
1180 *Journal of Communication* 14, 86–97. doi:10.1111/j.1460-2466.1964.
1181 tb02352.x.
- 1182 Jenkins, J.J., 1965. Mediation theory and grammatical behavior, in: Rosen-
1183 berg, S. (Ed.), *Directions in psycholinguistics*. MacMillan, New York, pp.
1184 66–96.
- 1185 Jenkins, J.J., 1968. The challenge to psychological theorists, in: Dixon,
1186 T.R., Horton, D.L. (Eds.), *Verbal behavior and general behavior theory*.
1187 Prentice-Hall, Inc., Englewood Cliffs, N.J., pp. 538–549.
- 1188 Jenkins, J.J., 1974. Remember that old theory of memory? well, forget it.
1189 *American Psychologist* 29, 785–795. doi:10.1037/h0037399.

- 1190 Jenkins, J.J., Palermo, D.S., 1964. Mediation processes and the acquisition
1191 of linguistic structure. *Monographs of the Society for Research in Child*
1192 *Development* 29, 141–169. doi:10.2307/1165762.
- 1193 Johns, B.T., Jamieson, R.K., Crump, M.J.C., Jones, M.N., Mewhort, D.J.K.,
1194 2016a. The combinatorial power of experience, in: Papafragou, A., Grod-
1195 ner, D., Mirman, D., Trueswell, J.C. (Eds.), *Proceedings of the 38th An-*
1196 *annual Meeting of the Cognitive Science Society*, Cognitive Science Society,
1197 Austin, TX. pp. 1325–1330.
- 1198 Johns, B.T., Jones, M.N., Mewhort, D.J.K., 2016b. Experience as a free pa-
1199 rameter in the cognitive modeling of language, in: Papafragou, A., Grod-
1200 ner, D., Mirman, D., Trueswell, J.C. (Eds.), *Proceedings of the 38th An-*
1201 *annual Meeting of the Cognitive Science Society*, Cognitive Science Society,
1202 Austin, TX. pp. 1325–1330.
- 1203 Johns, B.T., Taler, V., Pisoni, D.B., Farlow, M.R., Hake, A.M., Kareken,
1204 D.A., Unverzagt, F.W., Jones, M.N., 2013. Using cognitive models to
1205 investigate the temporal dynamics of semantic memory impairments in
1206 the development of Alzheimer’s disease, in: West, R., Stewart, T. (Eds.),
1207 *Proceedings of the 12th International Conference on Cognitive Modeling*.
1208 Carleton University, Ottawa, Canada, pp. 23–28. URL: [http://iccm-](http://iccm-conference.org/2013-proceedings/papers/0004/index.html)
1209 [conference.org/2013-proceedings/papers/0004/index.html](http://iccm-conference.org/2013-proceedings/papers/0004/index.html).
- 1210 Jones, M.N., Kintsch, W., Mewhort, D.J.K., 2006. High-dimensional se-
1211 mantic space accounts of priming. *Journal of Memory and Language* 55,
1212 534–552. doi:10.1016/j.jml.2006.07.003.
- 1213 Jones, M.N., Mewhort, D.J.K., 2007. Representing word meaning and order
1214 information in a composite holographic lexicon. *Psychological Review* 114,
1215 1–37. doi:10.1037/0033-295X.114.1.1.
- 1216 Kelly, M.A., Blostein, D., Mewhort, D.J.K., 2013. Encoding structure in
1217 holographic reduced representations. *Canadian Journal of Experimental*
1218 *Psychology* 67, 79–93. doi:10.1037/a0030301.
- 1219 Kelly, M.A., Kwok, K., West, R.L., 2015. Holographic declarative mem-
1220 ory and the fan effect: A test case for a new memory model for
1221 act-r, in: Taatgen, N.A., van Vugt, M.K., Borst, J.P., Mehlhorn,

- 1222 K. (Eds.), Proceedings of the 13th International Conference on Cog-
 1223 nitive Modeling. University of Groningen, Groningen, the Netherlands,
 1224 pp. 148–153. URL: [https://iccm-conference.neocities.org/2015/](https://iccm-conference.neocities.org/2015/proceedings/papers/0036/paper0036.pdf)
 1225 [proceedings/papers/0036/paper0036.pdf](https://iccm-conference.neocities.org/2015/proceedings/papers/0036/paper0036.pdf).
- 1226 Kelly, M.A., Reitter, D., 2017. Holographic declarative memory: Using
 1227 distributional semantics within act-r, in: Laird, J., Lebiere, C., Rosen-
 1228 bloom, P.S. (Eds.), The 2017 AAAI Fall Symposium Series: Technical
 1229 Reports, The AAAI Press, Palo Alto, California. pp. 382–387. URL:
 1230 <https://aaai.org/ocs/index.php/FSS/FSS17/paper/view/16001>.
- 1231 Kintsch, W., Van Dijk, T.A., 1978. Toward a model of text comprehension
 1232 and production. *Psychological Review* 85, 363–394. doi:10.1037/0033-
 1233 295X.85.5.363.
- 1234 Landauer, T.K., Dumais, S.T., 1997. A solution to plato’s problem: The la-
 1235 tent semantic analysis theory of acquisition, induction, and representation
 1236 of knowledge. *Psychological Review* 104, 211–240.
- 1237 Manning, C., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S., McClosky,
 1238 D., 2014. The Stanford CoreNLP natural language processing toolkit, in:
 1239 Proceedings of 52nd Annual Meeting of the Association for Computational
 1240 Linguistics: System Demonstrations, Association for Computational Lin-
 1241 guistics, Baltimore, Maryland. pp. 55–60. doi:10.3115/v1/P14-5010.
- 1242 McNeill, D., 1963. The origin of associations within the same grammatical
 1243 class. *Journal of Verbal Learning and Verbal Behavior* 2, 250 – 262. doi:10.
 1244 1016/S0022-5371(63)80091-2.
- 1245 Mewhort, D.J.K., Johns, B.T., Kelly, M., 2010. Applying the permutation
 1246 test to factorial designs. *Behavior Research Methods* 42, 366–372. doi:10.
 1247 3758/BRM.42.2.366.
- 1248 Mewhort, D.J.K., Kelly, M., Johns, B.T., 2009. Randomization tests and
 1249 the unequal-n/unequal-variance problem. *Behavior Research Methods* 41,
 1250 664–667. doi:10.3758/BRM.41.3.664.
- 1251 Mewhort, D.J.K., Shabahang, K.D., Franklin, D.R.J., 2018. Release from pi:
 1252 An analysis and a model. *Psychonomic bulletin & review* , 932–950doi:10.
 1253 3758/s13423-017-1327-3.

- 1254 Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J., 2013.
 1255 Distributed representations of words and phrases and their composi-
 1256 tionality, in: Burges, C.J.C., Bottou, L., Welling, M., Ghahramani,
 1257 Z., Weinberger, K.Q. (Eds.), Advances in Neural Information Pro-
 1258 cessing Systems 26. Curran Associates, Inc., pp. 3111–3119. URL:
 1259 [http://papers.nips.cc/paper/5021-distributed-representations-](http://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality.pdf)
 1260 [of-words-and-phrases-and-their-compositionality.pdf](http://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality.pdf).
- 1261 Murdock, B.B., 1982. A theory for the storage and retrieval of item and
 1262 associative information. *Psychological Review* 89, 609–626. doi:10.1037/
 1263 0033-295X.89.6.609.
- 1264 Nelson, K., 1977. The syntagmatic-paradigmatic shift revisited: A review of
 1265 research and theory. *Psychological Bulletin* 84, 93–116.
- 1266 Ororbia II, A.G., Mikolov, T., Reitter, D., 2017. Learning simpler language
 1267 models with the differential state framework. *Neural computation* 29,
 1268 3327–3352.
- 1269 Pennington, J., Socher, R., Manning, C.D., 2014. GloVe: Global vec-
 1270 tors for word representation, in: Empirical Methods in Natural Language
 1271 Processing (EMNLP), pp. 1532–1543. URL: [http://www.aclweb.org/](http://www.aclweb.org/anthology/D14-1162)
 1272 [anthology/D14-1162](http://www.aclweb.org/anthology/D14-1162).
- 1273 Pereira, F., 2000. Formal grammar and information theory: Together
 1274 again? *Philosophical Transactions of the Royal Society of London*
 1275 *A: Mathematical, Physical and Engineering Sciences* 358, 1239–1253.
 1276 doi:10.1098/rsta.2000.0583.
- 1277 Plate, T.A., 1995. Holographic reduced representations. *IEEE Transactions*
 1278 *on Neural Networks* 6, 623–641. doi:10.1109/72.377968.
- 1279 Pribram, K.H., 1969. The neurophysiology of remembering. *Scientific Amer-*
 1280 *ican* 220, 73–86.
- 1281 Princeton University, 2010. About wordnet. WordNet URL: [http://](http://wordnet.princeton.edu)
 1282 wordnet.princeton.edu.
- 1283 Reitter, D., Hockenmaier, J., Keller, F., 2006. Priming effects in combinatory
 1284 categorial grammar, in: Proceedings of the 2006 Conference on Empirical
 1285 Methods in Natural Language Processing, Association for Computational

1286 Linguistics, Stroudsburg, PA, USA. pp. 308–316. URL: <http://dl.acm.org/citation.cfm?id=1610075.1610119>.
1287

1288 Roberts, M.A., Chater, N., 2008. Using statistical smoothing to estimate the
1289 psycholinguistic acceptability of novel phrases. *Behavior Research Methods*
1290 40, 84–93.

1291 Rutledge-Taylor, M.F., Kelly, M.A., West, R.L., Pyke, A.A., 2014. Dy-
1292 namically structured holographic memory. *Biologically Inspired Cognitive*
1293 Architectures 9, 9–32. doi:10.1016/j.bica.2014.06.001.

1294 Saussure, F., 1916. Rapports syntagmatiques et rapports associatifs, in:
1295 Bally, C., Sechehaye, A. (Eds.), *Cours de linguistique générale*. Payot,
1296 Paris, France, pp. 170–175.

1297 Shannon, C.E., 1951. Prediction and entropy of printed english. *Bell System*
1298 Technical Journal 30, 50–64. doi:10.1002/j.1538-7305.1951.tb01366.
1299 x.

1300 Sidman, M., 2009. Equivalence relations and behavior: An introductory tuto-
1301 rial. *The Analysis of Verbal Behavior* 25, 5–17. doi:10.1007/bf03393066.

1302 Sloutsky, V.M., Yim, H., Yao, X., Dennis, S., 2017. An associative account
1303 of the development of word learning. *Cognitive Psychology* 97, 1 – 30.
1304 doi:10.1016/j.cogpsych.2017.06.001.

1305 Steedman, M., Baldridge, J., 2011. Combinatory categorial grammar, in:
1306 Borsley, R., Borjars, K. (Eds.), *Non-Transformational Syntax: Formal and*
1307 *Explicit Models of Grammar*. Wiley-Blackwell, pp. 181–224.

1308 Toutanova, K., Klein, D., Manning, C.D., Singer, Y., 2003. Feature-rich part-
1309 of-speech tagging with a cyclic dependency network, in: *Proceedings of*
1310 the 2003 Conference of the North American Chapter of the Association for
1311 Computational Linguistics on Human Language Technology, Association
1312 for Computational Linguistics, Edmonton, Canada. pp. 173–180. doi:10.
1313 3115/1073445.1073478.

1314 Ward, G., 1996. *Moby Part-of-Speech*. University of Sheffield. URL: <http://icon.shef.ac.uk/Moby/mpos.html>.
1315

1316 Appendix: Encoding order with one versus two permutations

1317 Our approach to encoding the sequential order of words differs from Jones
 1318 and Mewhort (2007). Convolution is commutative, that is, invariant to the
 1319 sequential order of the operands, i.e., $\mathbf{v}_1 * \mathbf{v}_2 = \mathbf{v}_2 * \mathbf{v}_1$. However, the se-
 1320 quence of the words can be preserved by permuting each operand. Jones and
 1321 Mewhort, p.35, using a method proposed by Plate (1995, p.12), apply two
 1322 different permutations to the left and right operands of convolution, such
 1323 that $(\mathbf{P}_{\text{left}}\mathbf{v}_1) * (\mathbf{P}_{\text{right}}\mathbf{v}_2) \neq (\mathbf{P}_{\text{left}}\mathbf{v}_2) * (\mathbf{P}_{\text{right}}\mathbf{v}_1)$.

1324 We apply a permutation only to the left operand as it is simpler and
 1325 sufficient for preserving sequence: $(\mathbf{P}_{\text{before}}\mathbf{v}_1) * \mathbf{v}_2 \neq (\mathbf{P}_{\text{before}}\mathbf{v}_2) * \mathbf{v}_1$. Our
 1326 one-permutation method is isomorphic to using two permutations. Vectors
 1327 constructed using one permutation will have, in expectation, the same spatial
 1328 relationships to each other as vectors constructed using two permutations,

$$\begin{aligned} &\text{cosine}((\mathbf{P}_{\text{before}}\mathbf{v}_1) * \mathbf{v}_2, (\mathbf{P}_{\text{before}}\mathbf{v}_3) * \mathbf{v}_4) \approx \\ &\text{cosine}((\mathbf{P}_{\text{right}}\mathbf{v}_1) * (\mathbf{P}_{\text{left}}\mathbf{v}_2), (\mathbf{P}_{\text{right}}\mathbf{v}_3) * (\mathbf{P}_{\text{left}}\mathbf{v}_4)) \end{aligned}$$

1329 where spatial relationships are measured by the cosine similarity. Differences
 1330 in the cosine similarity between the two methods will be due to small, zero
 1331 mean variations introduced by using the $\mathbf{P}_{\text{right}}$ permutation. The isomor-
 1332 phism arises because convolution and permutation preserve cosine similarity
 1333 relationships, such that $\text{cosine}(\mathbf{v}_1, \mathbf{v}_2) = \text{cosine}(\mathbf{P}\mathbf{v}_1, \mathbf{P}\mathbf{v}_2)$ and,

$$\text{cosine}((\mathbf{P}\mathbf{v}_1) * \mathbf{v}_2, (\mathbf{P}\mathbf{v}_3) * \mathbf{v}_4) \approx \text{cosine}(\mathbf{P}\mathbf{v}_1, \mathbf{P}\mathbf{v}_3) \times \text{cosine}(\mathbf{v}_2, \mathbf{v}_4) \quad (3)$$

1334 for any vectors $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4$ and permutation \mathbf{P} .